



LIBRARY OF THE  
UNIVERSITY OF ILLINOIS  
AT URBANA-CHAMPAIGN

510.84

I l 6r

no. 782-787

cop. 2



The person charging this material is responsible for its return to the library from which it was withdrawn on or before the **Latest Date** stamped below.

Theft, mutilation, and underlining of books are reasons for disciplinary action and may result in dismissal from the University.

UNIVERSITY OF ILLINOIS LIBRARY AT URBANA-CHAMPAIGN

JUN 6 1973

MAY 9 RECD



Digitized by the Internet Archive  
in 2013

<http://archive.org/details/statisticaltheor787park>

662  
0.787  
pa  
Report No. UIUCDCS-R-76-787

Math  
9  
NSF-OCA-DCR73-07980 A02-000018

THE STATISTICAL THEORY OF RELATIVE ERRORS  
IN FLOATING-POINT COMPUTATION

by

Douglass Stott Parker, Jr.

March 1976



DEPARTMENT OF COMPUTER SCIENCE  
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN · URBANA, ILLINOIS



Report No. UIUCDCS-R-76-787

THE STATISTICAL THEORY OF RELATIVE ERRORS  
IN FLOATING-POINT COMPUTATION

by

Douglass Stott Parker, Jr.

March 1976

Department of Computer Science  
University of Illinois at Urbana-Champaign  
Urbana, Illinois 61801

\* This work was supported in part by the National Science Foundation under Grant No. US NSF DCR73-07980 A02 and was submitted in partial fulfillment of the requirements for the degree of Master of Science in Computer Science, March 1976.





## TABLE OF CONTENTS

	Page
1. INTRODUCTION. . . . .	1
2. RELATIVE ERRORS AS PROBABILITY DENSITIES. . . . .	3
3. WHY ACCUMULATED ERROR DENSITIES ARE ALMOST NORMAL . . . .	9
4. PROOF THAT STATISTICAL BOUNDS HOLD FOR ROUND OFF ERRORS. .	20
5. CORRELATED ERRORS . . . . .	26
6. THE SNORT SYSTEM AND PROBLEMS OF IMPLEMENTATION . . . . .	28
7. CONCLUSIONS . . . . .	61
REFERENCES . . . . .	62



## 1. INTRODUCTION

For the man on the street, roundoff analysis of floating-point computation is synonymous with keeping bounds on accumulated errors: at any stage in some computation, the computed approximation  $f1(A)$  to a partial result  $A$  satisfies

$$f1(A) = A(1 + \delta)$$

where  $\delta = (f1(A) - A)/A$  is the relative error of the computation [16]. The idea is to find limits on the magnitude of  $\delta$ .

Recently, however, a number of papers have appeared adding probabilistic considerations to the distribution of roundoff errors in single arithmetic operations and of--what almost amounts to the same--errors in representing the reals as floating-point numbers ([3], [6], [14]). They show that relative errors are not uniformly distributed within their bounds, but instead cluster around zero. The error bounds themselves are very pessimistic.

In this paper relative errors  $\delta$  are viewed as random variables with suitable probability distributions. The results of applying classical probability theory to roundoff analysis in this fashion are: (1) A proof of the generally accepted statement that accumulated errors approach normal distributions after several operations ([13, p. 113], [11, p. 104], [4, p. 306]). (2) An appreciation of the rate of convergence to normal distributions and how floating-point addition and correlated errors can slow this convergence down. (3) A new method of

automated error analysis comparable to (but in most ways inferior to) Interval Analysis, but in every way superior to Wilkinson bound analysis.

We should comment that it is not "unrealistic" to treat relative errors as random variables. To do so is just to follow the basic philosophy of Bayesian statistics: if you don't know the exact value of a variable, but know what its possible values are and what the probability of taking each of those values on is, then you should treat the variable as random with the corresponding probability distribution. Treating roundoff errors in this way makes our work fairly clean, plus lets it be general and rigorous as well.

## 2. RELATIVE ERRORS AS PROBABILITY DENSITIES

In the analysis of Wilkinson [16] the relative error

$\delta_A = (f1(A) - A)/A$  in a computed result  $A$  is rarely known exactly-- usually one says that it lies within known bounds. These bounds are proportional to the machine precision  $\beta^{-t}$ , where  $\beta$  is the arithmetic base and  $t$  is the number of mantissa digits used in representing floating-point numbers. Computation of bounds is simple: if the floating-point approximations to  $A$  and  $B$  involve the relative errors  $\delta_A$  and  $\delta_B$ , respectively, then the floating-point approximation to  $A \cdot B$  is

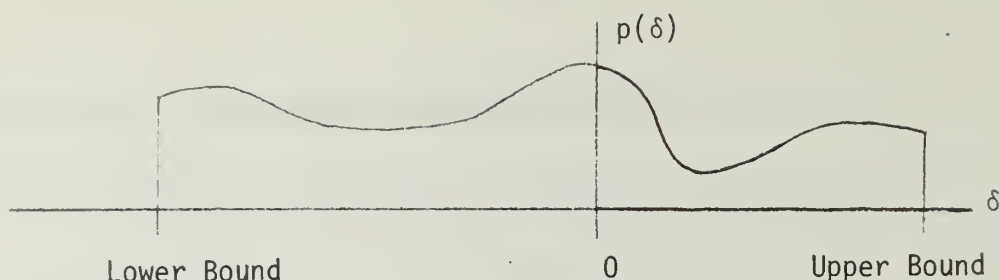
$$\begin{aligned} f1(A \cdot B) &= [A(1 + \delta_A)][B(1 + \delta_B)](1 + \delta_R) \\ &= AB(1 + \delta_A + \delta_B + \delta_R + \delta_A\delta_B + \delta_A\delta_R \\ &\quad + \delta_B\delta_R + \delta_A\delta_B\delta_R) \end{aligned} \tag{2.1}$$

where  $\delta_R$  is the rounding error incurred in the multiplication and is bounded by the machine precision. Since hopefully  $|\delta_A| \ll 1$ ,  $|\delta_B| \ll 1$ , we have

$$f1(AB) \approx AB(1 + \delta_A + \delta_B + \delta_R) \tag{2.2}$$

to a very good approximation. Then the bound for the error on  $f1(AB)$  is the sum of the bounds for  $\delta_A$ ,  $\delta_B$ , and  $\delta_R$ .

Up to this point we have treated each relative error  $\delta$  as an unknown element of a known interval. It is not a big step to view  $\delta$  as being smeared over this interval, having at each point a probability of occurrence:



Thus we equate  $\delta$  with an underlying probability density  $p(\delta)$  so that

$$\begin{aligned} \text{Probability } (\delta \leq x_0) &= \int_{-\infty}^{x_0} p(x) dx \\ &= \int_{\text{Lower Bound}}^{x_0} p(x) dx \end{aligned} \quad (2.3)$$

Recall that a probability density  $f$  is a function on  $\mathbb{R}$  satisfying

$$\begin{aligned} \text{(i)} \quad & f(x) \geq 0 \quad \forall x \in \mathbb{R} \\ \text{(ii)} \quad & \int_{-\infty}^{\infty} f(x) dx = 1 \end{aligned} \quad (2.4)$$

and that the probability distribution  $F$  corresponding to the density  $f$  is defined by

$$F(x) = \int_{-\infty}^x f(t) dt \quad (2.5)$$

The question arises as to what expressions like " $A(1 + \delta_A)$ " mean if we view  $\delta_A$  as a probability density. This is resolved by the following definition and explanation.

Def: A smear number (or fuzzy number) is a pair  $N = [A, f]$  where  $A \in \mathbb{R}$  and  $f$  is a probability density on  $\mathbb{R}$ . The probability that  $N$  takes on a value  $\leq A \cdot C$  for any  $C \in \mathbb{R}$  is given by

$$\Pr(N \leq A \cdot C) = \int_{-\infty}^C f(t) dt \quad (2.6)$$

The idea is that every floating-point approximation to a number is a smeared number. Thus  $fl(A) = A(1 + \delta_A) = [A, (1 + \delta_A)]$ .

We clarify this concept with three observations:

Observation 1. The constant 1 in expressions of the form  $(1 + \delta)$  should be thought of as an impulse function  $\Delta$  (Dirac delta) at 1.

This is natural since 1 should represent a density with all its support at 1. In the following we will use 1 and  $\Delta$  interchangeably when no confusion should result.

Observation 2. If  $\delta_1$  and  $\delta_2$  are independent random variables with corresponding densities  $p_1$  and  $p_2$ , then  $\delta = \delta_1 + \delta_2$  corresponds to the convolution density  $p = p_1 * p_2$ .

Since  $\delta$  is the density which gives the probability at each point  $x$  that  $\delta_1 + \delta_2 = x$ ,

$$\begin{aligned} \Pr(\delta \leq x) &= \iint_{y+z \leq x} p_1(y) p_2(z) dydz \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{x-z} p_1(y) p_2(z) dydz \end{aligned}$$

Differentiating with respect to  $x$ ,

$$p(x) = \int_{-\infty}^{\infty} p_1(x-z) p_2(z) dz = p_1 * p_2(x)$$

Note from observations 1 and 2 that if  $p_A$  is the density corresponding to  $\delta_A$ , then

$$f1(A) = A(1 + \delta_A) = [A, (1 + \delta_A)] = [A, \Delta * p_A]$$

Observation 3. If  $\delta_1$  and  $\delta_2$  are independent random variables with corresponding densities  $p_1$  and  $p_2$ , then  $\delta = \delta_1 \delta_2$  corresponds to the density  $p$  defined by

$$p(x) = \int_{-\infty}^{\infty} p_1(x/z) p_2(z) \frac{dz}{z}.$$

The derivation is similar to that for Observation 2. We remark that the approximation

$$(1 + \delta_1 + \delta_2 + \delta_1 \delta_2) \approx (1 + \delta_1 + \delta_2)$$

when  $\delta_1$  and  $\delta_2$  are small random variables (with corresponding densities) still makes sense: if  $p_1$  is zero outside  $[-D_1, D_1]$  and  $p_2$  is zero outside  $[-D_2, D_2]$  then the density  $p$  defined as above for  $\delta_1 \delta_2$  is zero outside  $[-D_1 D_2, D_1 D_2]$ . Assuming  $D_1, D_2 \ll 1$ ,  $p$  approximates an impulse function at zero (the identity for convolution).

We now consider what  $f1(A + B)$  looks like. Since

$$\begin{aligned} f1(A + B) &= (A(1 + \delta_A) + B(1 + \delta_B)) (1 + \delta_R) \\ &= A(1 + \delta_A + \delta_R) + B(1 + \delta_A + \delta_R) \end{aligned}$$

where again  $\delta_R$  is rounding error, we start with

$$f1(A + B) \approx [A, 1 * p_A * p_R] + [B, 1 * p_B * p_R]$$

The results for subtraction and division are similar. We summarize the results here:



$$f1(A + B) = [A, 1 * p_A * p_R] + [B, 1 * p_B * p_R] \quad (2.7)$$

$$f1(A - B) = [A, 1 * p_A * p_R] + [-B, 1 * p_B * p_R] \quad (2.8)$$

$$f1(A \cdot B) = [A \cdot B, 1 * p_A * p_B * p_R] \quad (2.9)$$

$$f1(A/B) = [A/B, 1 * p_A * p_B^V * p_R] \quad (2.10)$$

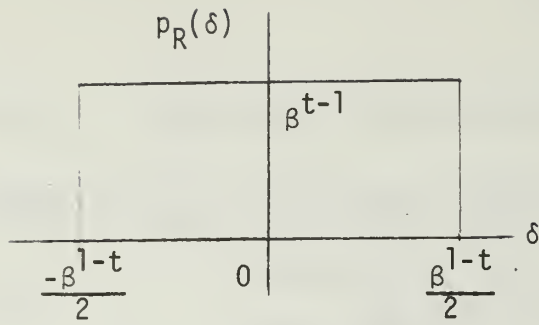
In (2.10),  $p_B^V$  is defined by  $p_B^V(x) = p_B(-x)$ . Notice that in (2.7) and (2.8) the result is not actually a smeared number but a sum of smeared numbers. This reflects Wilkinson's statement [16, §28] that you cannot say much about the error in a sum without knowing more about the magnitudes of the summands. In the next section we show a method to reduce the sum to a single smeared number. In short, we show that we can make a consistent algebra for smeared numbers.

We must now select probability densities  $p(\delta)$  for (1) representation errors (i.e., errors incurred by representing real numbers in floating-point format) and (2) the rounding error  $\delta_R$  made in single arithmetic operations. All other errors are sums of these two kinds. As mentioned in the introduction, we can use the same density for both (1) and (2):

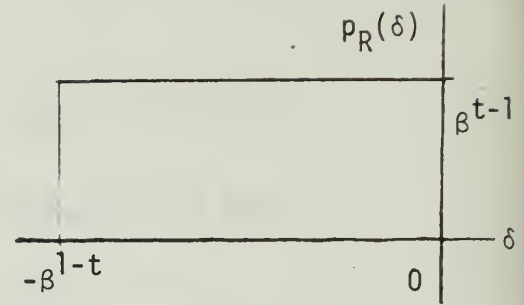
Let  $p_R$  be the uniform probability density on

$$[-\beta^{1-t}/2, \beta^{1-t}/2] \text{ if } \underline{\text{rounded}} \text{ arithmetic is used} \quad (2.11)$$

$$[-\beta^{1-t}, 0] \text{ if } \underline{\text{chopped}} \text{ (truncated) arithmetic is used}$$



$p_R$  for Rounded Arithmetic



$p_R$  for Chopped Arithmetic

Note that  $p_R$  is a uniform distribution inside the ordinary Wilkinsonian bounds. This is very pessimistic but is the best one can do a priori.

If we consider the error in representing  $A = \alpha \cdot \beta^e$  where  $\alpha \in [1/\beta, 1)$ ,  $e \in \mathbb{Z}$  then since

$$[fl(\alpha) - \alpha] \quad (2.12)$$

is uniformly distributed on  $[-\beta^{-t}/2, \beta^{-t}/2]$  for rounding or  $[-\beta^{-t}, 0]$  for chopping (because the "(t + 1)-th digits" is practically uniform; cf. [14, p. 270])

$$\delta_A = (fl(A) - A)/A = (fl(\alpha) - \alpha)/\alpha \quad (2.13)$$

is uniformly distributed on  $[-\beta^{-t}/2\alpha, \beta^{-t}/2\alpha]$  for rounding or  $[-\beta^{-t}/\alpha, 0]$  for chopping. If we assume the worst case  $\alpha = 1/\beta$ , we get the  $p_R$  defined above.

### 3. WHY ACCUMULATED ERROR DENSITIES ARE ALMOST NORMAL

This section introduces the Central Limit Theorem, which is the basis for the popular statement that accumulated relative errors become normally distributed. The Central Limit Theorem is useful intuitively but we avoid its proof, offering Section 4 instead. In this section we explain how a sum of relative errors converges to a normally distributed error. After showing that sums of smeared numbers are smeared numbers we get an appreciation of how this convergence may be disturbed.

We begin by reviewing some important concepts from probability theory. Let  $\xi$  be a random variable having corresponding probability density,  $f$ . Then:

Def: The nth moment  $M_n$  of  $\xi$  (or of  $f$ ) is

$$M_n = \int_{-\infty}^{\infty} x^n f(x) dx. \quad (3.1)$$

Def: Let  $g$  be an arbitrary complex-valued function on  $\mathbb{R}$ . Define the expected value of  $g(\xi)$  by

$$E[g(\xi)] = \int_{-\infty}^{\infty} g(x)f(x) dx \quad (3.2)$$

Note therefore that  $M_n = E[\xi^n]$ .

Def: The mean of  $\xi$  (or of  $f$ ) is its first moment:

$$\mu = E[\xi] = M_1 \quad (3.3)$$

Def: The variance of  $\xi$  (or of  $f$ ) is given by

$$\sigma^2 = E[(\xi - \mu)^2] = M_2 - M_1^2 \quad (3.4)$$

The definitions (3.1) - (3.4) are tied together with the following theorem:

Theorem 1. Given densities  $f_1, \dots, f_N$  where  $f_i$  has mean  $\mu_i$  and variance  $\sigma_i^2$ , then the convolution density  $f_1 * \dots * f_N$  has mean  $\mu = \sum_{i=1}^N \mu_i$  and variance  $\sigma^2 = \sum_{i=1}^N \sigma_i^2$ .

Proof. Compute the Fourier transforms ("characteristic functions,"  $E[\exp(iw\xi)]$ ) of both sides of  $f = f_1 * \dots * f_N$  and use the facts that

$$(1) \quad \hat{g}(w) = \sum_{n=0}^{\infty} \frac{(-i)^n M_n}{n!} \quad (i = \sqrt{-1}; \hat{\phantom{x}} = \text{Fourier Trans- form})$$

$$(2) \quad g * h(w) = \hat{g}(w) \hat{h}(w)$$

Now recall that the Normal (or Gaussian) probability density with parameters  $\mu, \sigma (\sigma > 0)$  is

$$\psi(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-(x - \mu)^2 / 2\sigma^2} \quad (3.5)$$

so its probability distribution  $\Phi(x)$  is

$$\Phi(x) = \frac{1}{\sqrt{2\pi} \sigma} \int_{-\infty}^x e^{-(t - \mu)^2 / 2\sigma^2} dt \quad (3.6)$$

It is easy to show that the mean and variance of  $\psi$  are  $\mu$  and  $\sigma^2$ , respectively.

The Central Limit Theorem may be stated as follows: given a set  $\{\delta_1, \dots, \delta_N\}$  of random variables having respective probability densities  $p_i$ , means  $\mu_i$ , and variances  $\sigma_i^2$  ( $i = 1, \dots, N$ ) then under "certain loose conditions" the sum variable  $\delta = \delta_1 + \delta_2 + \dots + \delta_N$  (having density  $p = p_1 * p_2 * \dots * p_N$ ) converges as  $N \rightarrow \infty$  to a normal density with mean  $\mu = \mu_1 + \mu_2 + \dots + \mu_N$  and variance  $\sigma^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_N^2$ .

Proof of the Central Limit Theorem is dependent of course on the "loose conditions" surrounding the summand variables and the type of convergence used. Good surveys may be found in [1] and [12]. We will omit any proof here for roundoff errors because the convergence to a normal distribution, although useful intuitively, is difficult to prove and unnecessary for obtaining statistical bounds (see Section 4).

The idea is that the smeared number  $A(1 + \delta_1 + \delta_2 + \dots + \delta_N)$  is equal to  $A(1 + \delta)$  where  $\delta$  has a density that is almost normal.

Since:

1. We can compute  $\delta$ 's mean  $\mu$  and variance  $\sigma^2$  easily.
2. Assuming  $\delta$  normal, the probability that  $|\delta - \mu| = 3\sigma$  is 0.002798 (i.e., a 99.7% confidence interval on  $\delta$ 's value is  $[\mu - 3\sigma, \mu + 3\sigma]$ ).
3. For  $N$  larger than 2 or 3 the error bound  $[\mu - 3\sigma, \mu + 3\sigma]$  is inside the crude Wilkinsonian error bound  $[-N\beta^{1-t}/2, N\beta^{1-t}/2]$  for Rounding or  $[-N\beta^{1-t}, 0]$  for truncation

--apparently without much effort we can get much tighter, more realistic error bounds and improve our error estimates.

This gives us our motivation: statistical error bounds seem tight. But we need more details.

We compute the mean and variance for the rounding error density  $p_R$  defined by (2.11). For rounding, we have

$$\begin{aligned}\mu_{\text{Rounding}} &= 0 \\ \sigma_{\text{Rounding}}^2 &= (\beta^{1-t}/2)^2/3 = \beta^{2-2t}/12\end{aligned}\quad (3.7)$$

For chopping, similarly

$$\begin{aligned}\mu_{\text{Chopping}} &= -\beta^{1-t}/2 \\ \sigma_{\text{Chopping}}^2 &= \frac{(\beta^{1-t})^2}{3} - \frac{(-\beta^{1-t})^2}{2} = \beta^{2-2t}/12\end{aligned}\quad (3.8)$$

Therefore

$$\sigma_R \stackrel{\text{def}}{=} \sigma_{\text{Rounding}} = \sigma_{\text{Chopping}} = \beta^{1-t}/\sqrt{12}\quad (3.9)$$

Notice that from (2.9) and (2.10) expressions involving only repeated multiplication, repeated division, or mixed multiplication and division will result in smeared numbers of the form  $[A, 1 * \overbrace{p_R * p_R * p_R * \dots p_R}^m]$  where the number  $m$  of  $p_R$ 's in the convolution product depends of course on the number of operands in the expression and the accumulated error for each operand. For example if  $A, B, C$  are irrational then  $\text{fl}(A(B \cdot C)) = [ABC, 1 * p_R * (p_R * p_R * p_R) * p_R]$  since each number has a representation error bounded by  $p_R$ , and the two multiplication operations have



rounding (truncation) errors bounded by  $p_R$ . The Wilkinson bounds for the relative error in ABC would be  $\delta_{ABC} \in [-\frac{5}{2} \beta^{1-t}, \frac{5}{2} \beta^{1-t}]$  for rounding and  $\delta_{ABC} \in [-5\beta^{1-t}, 0]$  for truncation.

If we assume (Central Limit Theorem) that  $(p_R * p_R * p_R * p_R * p_R)$  is normal with  $\sigma^2 = 5\sigma_R^2$ ,  $\mu = 0$  for rounding or  $\mu = 5(-1/2\beta^{1-t})$  for chopping, then the error in ABC should satisfy

$$\delta_{ABC} \in [\mu - 3\sigma, \mu + 3\sigma] = \begin{cases} [-\sqrt{\frac{5}{12}} \beta^{1-t}, 3\sqrt{\frac{5}{12}} \beta^{1-t}] & \text{Rounding} \\ [(\frac{-5}{12} - 3\sqrt{\frac{5}{12}})\beta^{1-t}, (\frac{-5}{2} + 3\sqrt{\frac{5}{12}})\beta^{1-t}] & \text{Chopping} \end{cases}$$

99.7% of the time. Now  $3\sqrt{\frac{5}{12}} \approx 1.94 < 5/2$ , and we have improved our error bounds by 20% after only two operations. As the number of operations  $N$  increases we will get better improvement since Wilkinson's bound are proportional to  $N$  while the  $3\sigma$  bounds are proportional to  $\sqrt{N}$ .

It is not unreasonable to assume that  $(p_R * p_R * p_R * p_R * p_R)$  is normal, either. Figures 1-4 show the result of repeated convolution with the error density for Rounding. Already  $(p_R * p_R * p_R * p_R)$  is essentially normal (Figure 4). On a superficial basis at the least, convergence to a normal density seems rapid.

It should be pointed out that the use of division introduces "checked" densities  $p_R^V$  in the convolution product (c.f. (2.10)). This is not a problem because

$$p_R^V(x) = p_R(-x) = \begin{cases} p_R(x) & \text{for Rounding} \\ p_R(x - \beta^{1-t}) & \text{for Chopping} \end{cases} \quad (3.10)$$

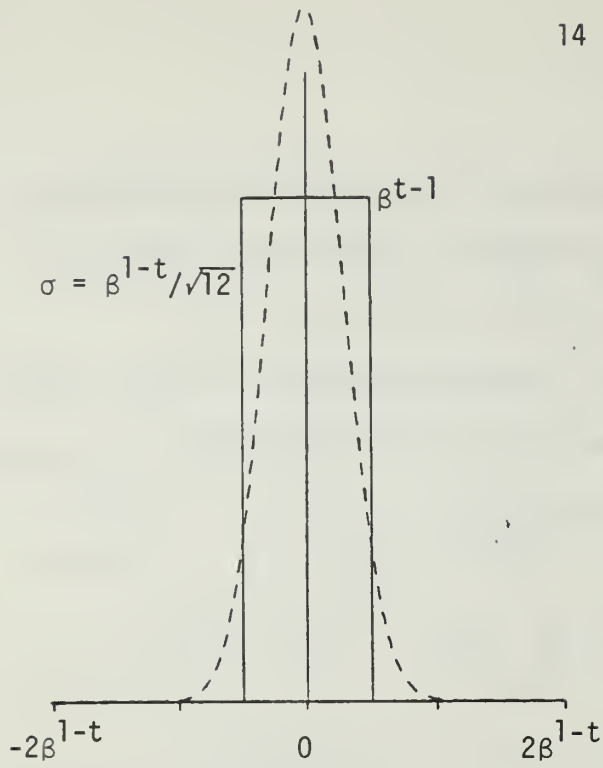


Figure 1.

$p_R$  for Rounding and its normal approximation

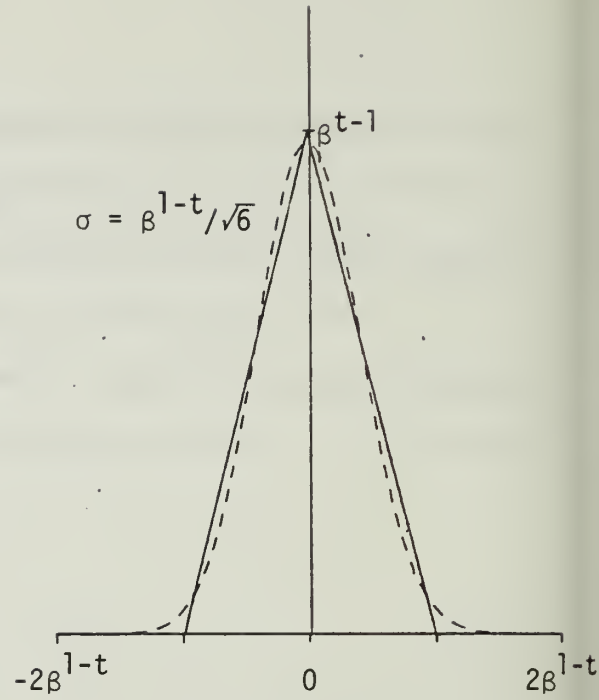


Figure 2.

$p_R * p_R$  and normal approximation

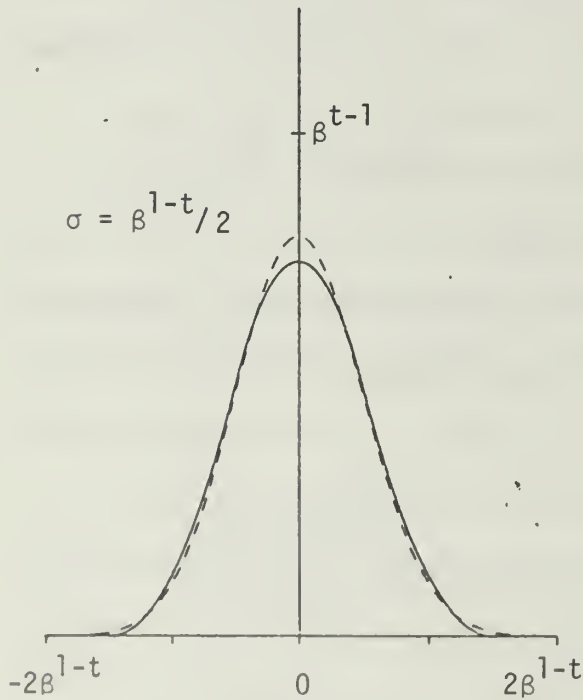


Figure 3.

$p_R * p_R * p_R$  and normal approximation

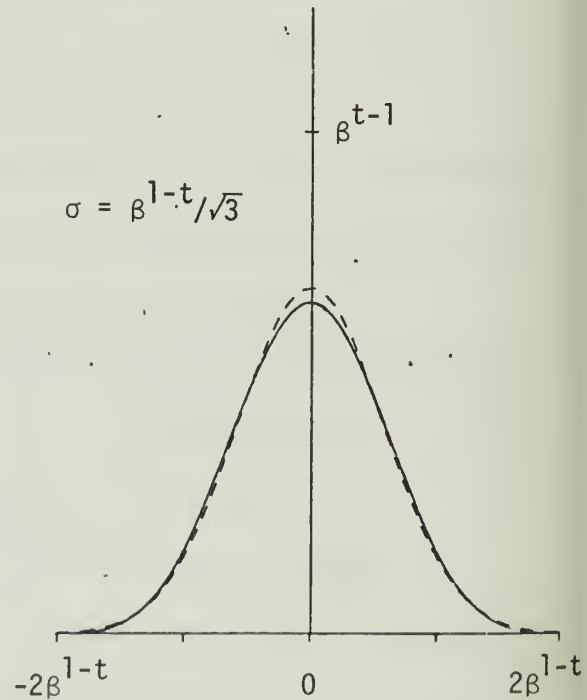


Figure 4.

$p_R * p_R * p_R * p_R$  and normal approximation



Obviously what was said in the previous paragraph holds without modification for Rounding. For chopping we must simply be more careful in computation of the mean  $\mu$ , since the mean of  $p_R^v$  is  $+\frac{1}{2}\beta^{1-t}$  instead of  $-\frac{1}{2}\beta^{1-t}$  as it was before. For Chopping, division errors tend to cancel out multiplication errors since they are of the opposite sign--this is reflected here by the fact that division makes the mean of the final error more positive, while multiplication makes it more negative.

We commented in Section 2 that addition and subtraction were not as nice as multiplication and division insofar as their results were not easily expressible as smeared numbers. For the rest of this section we address the problem of reducing a sum of smeared numbers to a single smeared number.

It is expedient to work in absolute errors instead of the relative errors we have been working with until now. The following definition and lemma provide the necessary groundwork.

Def: For a bounded density  $f$  and nonzero real number  $A$ , the stretched density  $f^A$  is given by

$$f^A(x) = \frac{1}{|A|} f\left(\frac{x}{A}\right) \quad (3.11)$$

We say that  $f^A$  is  $f$  stretched by the value  $A$ . Here we have restricted our definition to bounded densities to avoid the problem introduced by impulse functions.

Lemma 1 If  $f$  is a bounded density, then  $[A, f] = [1, f^A]$ .

Proof Since  $\int_{-\infty}^{A \cdot C} f^A(t) dt = \int_{-\infty}^{A \cdot C} \frac{1}{|A|} f\left(\frac{t}{A}\right) dt \Big|_{\mu=t/A} = \int_{-\infty}^C f(\mu) d\mu.$

$[1, f^A]$  is equivalent to  $[A, f]$  under the definition of smeared numbers. Notice that  $[1, f^A]$  is an "absolute" smeared number, no longer relative to  $A$ . Effectively  $[1, f^A] = f^A$ .

Theorem 2 If  $f1(A + B) = [A, 1 * p_1] + [B, 1 * p_2]$  and  $A \neq 0$ ,

$B \neq 0$ ,  $A + B \neq 0$  then  $f1(A + B) = [A + B, 1 * p]$  where  $p(x)$

$$= p_1^{(A/A+B)} * p_2^{(B/A+B)}(x) \quad (3.12)$$

Proof Brute force. Note that  $(1 * f)(x) = f(x - 1)$  and  $(1 * f)^A(x) = \frac{1}{|A|} f\left(\frac{x}{A} - 1\right)$ . Now:

$$\begin{aligned} f1(A + B) &= [A, 1 * p_1] + [B, 1 * p_2] \\ &= [1, (1 * p_1)^A] + [1, (1 * p_2)^B] \\ &= (1 * p_1)^A * (1 * p_2)^B \quad (\text{Lemma 1}) \end{aligned}$$

(This follows from the remark that  $[1, f] = f$ , and observation 2, Part 2.)

$$\begin{aligned} &= \int_{-\infty}^{\infty} \frac{1}{|A|} p_1\left(\frac{x-t}{A} - 1\right) \frac{1}{|B|} p_2\left(\frac{t}{B} - 1\right) dt \Big|_{t=(A+B)\mu} \\ &= \left| \frac{A+B}{AB} \right| \int_{-\infty}^{\infty} p_1\left(\frac{x - (A+B)\mu - A}{A}\right) \\ &\quad p_2\left(\frac{(A+B)\mu - B}{B}\right) d\mu \Big|_{\mu=v+(B/A+B)} \\ &= \left| \frac{A+B}{AB} \right| \int_{-\infty}^{\infty} p_1\left(\frac{x - (A+B)v - (A+B)}{A}\right) \\ &\quad p_2\left(\frac{(A+B)v}{B}\right) dv \end{aligned}$$

$$\begin{aligned}
&= \left| \frac{A+B}{AB} \right| \int_{-\infty}^{\infty} \left| \frac{A}{A+B} \right| p_1^{(A/A+B)} \left( \frac{x}{A+B} - 1 - v \right) \\
&\quad \left| \frac{B}{A+B} \right| p_2^{(B/A+B)}(v) dv \\
&= \frac{1}{|A+B|} p_1^{(A/A+B)} * p_2^{(B/A+B)} \left( \frac{x}{A+B} - 1 \right) \\
&= \frac{1}{|A+B|} 1 * p_1^{(A/A+B)} * p_2^{(B/A+B)} \left( \frac{x}{A+B} \right) \\
&= (1 * p_1^{(A/A+B)} * p_2^{(B/A+B)})^{A+B} \\
&= [1, (1 * p_1^{(A/A+B)} * p_2^{(B/A+B)})^{A+B}] \\
&= [A+B, 1 * p_1^{(A/A+B)} * p_2^{(B/A+B)}]
\end{aligned}$$

Notice that  $A - B = A + (-B)$  so Theorem 2 applies to subtraction also. Theorem 2 leaves us a few steps short of finishing off all the problems of addition.

### Corollary 1.

$$\begin{aligned}
\text{If} \quad & f1\left(\sum_{i=1}^n A_i\right) = [A_1, 1 * p_1] + [A_2, 1 * p_2] + \dots + [A_n, 1 * p_n] \\
\text{then} \quad & f1\left(\sum_{i=1}^n A_i\right) = \left[ \sum_{i=1}^n A_i, 1 * p \right] \\
\text{where} \quad & p = p_1^{(A_1/\sum A_i)} * p_2^{(A_2/\sum A_i)} * \dots * p_n^{(A_n/\sum A_i)} \quad (3.13)
\end{aligned}$$

(Proof. Trivial. Naturally we now require that  $A_i \neq 0$  ( $i = 1, \dots, n$ ) and  $\sum_{i=1}^n A_i \neq 0$ .)

Theorem 3. Let the mean and variance for the density  $f$  be  $\mu$  and  $\sigma^2$ , respectively. Then the mean and variance for  $f^A$  are  $(A\mu)$  and  $(A^2\sigma^2)$ , respectively. (Proof. Trivial.)

From corollary 1 and Theorem 3 it is immediately obvious how to apply the central limit theorem to addition and subtraction. If the relative error density  $p_1$  of  $A$  has mean  $\mu_1$  and variance  $\sigma_1^2$  and the density  $p_2$  of  $B$  has mean  $\mu_2$  and variance  $\sigma_2^2$ , then the density  $p$  of  $A + B$  (as in Theorem 1) has

$$\text{mean} \quad \mu = \left(\frac{A}{A+B}\right)\mu_1 + \left(\frac{B}{A+B}\right)\mu_2 \quad (3.14)$$

$$\text{and variance } \sigma^2 = \left(\frac{A}{A+B}\right)^2\sigma_1^2 + \left(\frac{B}{A+B}\right)^2\sigma_2^2 \quad (3.15)$$

After many additions the error density will converge to a normal probability density--but more operations are needed because the stretched densities  $p_j^{(A_j/\sum A_i)}$  will not be identical up to translation as they were for multiplication: accumulated errors in multiplication become normally distributed much more quickly than those in addition. The reader who has researched the Central Limit Theorem has probably discovered that most of the existing proofs assume identically distributed summand variables. Unidentical distributions can slow down convergence--especially when the variables differ greatly. Thus when stretching factors like  $A/(A - B)$  became large (implying lots of cancellation and loss of significance in the computation of  $A - B$ ), convergence toward normal error densities is slowed.

In the next section we show that we can guarantee confidence intervals almost as good as  $[\mu - 3\sigma, \mu + 3\sigma]$  all of the time, even with

slowed convergence. Thus convergence and the Central Limit Theorem in general should be used only for an intuitive feel of what is going on, while the growth of  $\sigma$  is worth serious attention. Note that when stretching factors like  $A/(A - B)$  get large, convergence not only slows down but variances get large also (Theorem 3). This is a serious weakness of statistical error analysis.

## 4. PROOF THAT STATISTICAL BOUNDS HOLD FOR ROUND OFF ERRORS

Section 3 introduced the Central Limit Theorem, showed that it could be used to show how accumulated roundoff errors have asymptotically normal distributions, and discussed the rate of convergence to normal distributions. The literature surrounding the Central Limit Theorem is labyrinthine since there are any number of assumptions one can make about the summand variables which simplify or aggravate the proof. However, as we mentioned earlier we are not so much concerned with proving convergence as being able to say, for any roundoff variable  $\delta$ , something like

$$\Pr(|\delta - \mu| \geq 3\sigma) \leq .0028 \quad (4.1)$$

as we can with normally distributed random variables. In general given a fixed small probability  $\rho$  we want to find a corresponding  $\lambda$  such that

$$\Pr(|\delta - \mu| = \lambda\sigma) \leq \rho. \quad (4.2)$$

Since  $\delta$  will correspond to a lengthy convolution of error densities, determining  $\lambda$  exactly will be unfeasible. Instead we can get bounds on  $\lambda$ . It would be simpler to just assume that  $\delta$  is normally distributed and look up  $\lambda$  in normal distribution tables. However, this is not always legitimate: see e.g. [5]. Although three or four convolutions can approximate a normal density well (see Figures 1-4) large stretching factors can upset this approximation. The purpose of this section is to bound how much these factors can upset it, by obtaining bounds for  $\lambda$  in (4.2). We start with Chebyshev's Inequality.

Theorem 4. Let  $\delta$  be a random variable with underlying density  $p(t)$ , having mean 0 and variance  $\sigma^2$ . Then

$$\Pr(|\delta| \geq \lambda\sigma) \leq 1/\lambda^2$$

Proof.  $\sigma^2 = E[\delta^2]$

$$\begin{aligned} &= \int_{-\infty}^{\infty} t^2 p(t) dt \\ &\geq \int_{|t| \geq \lambda\sigma} t^2 p(t) dt \\ &\geq (\lambda\sigma)^2 \int_{|t| \geq \lambda\sigma} p(t) dt \\ &= (\lambda\sigma)^2 \Pr(|\delta| \geq \lambda\sigma) \end{aligned}$$

$$\Rightarrow \Pr(|\delta| \geq \lambda\sigma) \leq \frac{1}{\lambda^2}, \text{ as desired.}$$

This says that if we want

$$\Pr(|\delta| \geq \lambda\sigma) \leq .0028$$

then we will be guaranteed this relationship provided

$$\frac{1}{\lambda^2} = .0028 \Rightarrow \lambda \geq 18.89$$

Clearly this is not very good since  $\lambda = 3$  works for normally distributed variables, and after many convolutions we expect  $\delta$  to be very close to normal. It turns out we can obtain a much tighter limit on  $\lambda$  if we go through the analysis used to obtain the Chernoff bound.

Theorem 5. Let  $\delta = \sum_{i=1}^N \delta_i$  be a roundoff error, i.e.,  $\delta$  is represented by



the convolution of  $N$  (possibly stretched) copies of  $p_R$ , the rounding density. Then

$$\Pr(|\delta| \geq \lambda\sigma) \leq 2e^{-\lambda^2/2}$$

where  $\sigma^2$  is the variance of  $\delta$ .

Proof. Set  $d(x) = \begin{cases} 0 & x \leq \lambda\sigma \\ 1 & x \geq \lambda\sigma \end{cases}$

Then

$$d(x) \leq e(x) \stackrel{\text{def}}{=} \exp(L(x - \lambda\sigma)) \text{ for any positive } L.$$

This implies

$$\begin{aligned} \Pr(\delta \geq \lambda\sigma) &= E[d(\delta)] \\ &\leq E[e(\delta)] \\ &= E[\exp(L\delta) \exp(-L\lambda\sigma)] \\ &= E[\exp(L\delta)] e^{-L\lambda\sigma} \end{aligned}$$

Now if  $p(t)$  is the density associated with  $\delta$  then

$$\begin{aligned} E[\exp(L\delta)] &= \int_{-\infty}^{\infty} \exp(Lt) p(t) dt \\ &= \sum_{n=0}^{\infty} \frac{L^n}{n!} \int_{-\infty}^{\infty} t^n p(t) dt \quad [p \text{ is of compact support}] \\ &= \sum_{n=0}^{\infty} \frac{M_n L^n}{n!} \end{aligned}$$

where  $M_n$  is the  $n$ th moment of  $p(t)$ . Since  $p(t)$  is symmetric about 0 (i.e., even) the odd terms vanish and we are left with



$$E[\exp(L\delta)] = \sum_{n=0}^{\infty} \frac{M_{2n} L^{2n}}{(2n)!}$$

We now need the following lemma.

Lemma 2. For roundoff densities  $p(t)$

$$\frac{M_{2n}}{(M_2)^n} \leq \frac{(2n)!}{n! 2^n}$$

Proof. It can be shown that if  $p(t) = p_1 * p_2 * \dots * p_N(t)$  and  $M_{n,k}$  is the  $n$ th moment of  $p_k$ , that the  $n$ th moment  $M_n$  of  $p$  is

$$M_n = \sum_{\substack{n_1 + \dots + n_N = n}} \frac{n!}{n_1! n_2! \dots n_N!} \prod_{k=1}^N M_{n_k, k}$$

This is done as in Theorem 1 by comparing the characteristic functions (essentially Fourier transforms) of both sides of the equation

$p = \prod_{k=1}^N p_k$ . We then use the fact that the  $p_k$  are stretched versions of the roundoff density  $p_R$  (2.11), with moments

$$M_{n,k} = \begin{cases} \frac{[A_k(\beta^{1-t}/2)]^n}{(2n+1)} & n \text{ even} \\ 0 & n \text{ odd} \end{cases}$$

where  $p_k = p_R^{(A_k)}$ , i.e.,  $A_k$  is a nonzero stretching factor. Pulling all of this material together leads to the bound of the lemma.

It is interesting that if  $p(t)$  is normal with mean 0,

$$\frac{M_{2n}}{(M_2)^n} = \frac{M_{2n}}{\sigma^{2n}} = \frac{(2n)!}{n! 2^n} \quad (M_2 = \sigma^2)$$

--so we have equality with the stated bound. Therefore:

1. The bound is fairly tight, especially as N gets large.
2. The lemma holds for other densities besides stretched versions of  $p_R$ .
3. There is probably an elegant proof of this lemma.

It now follows that

$$E[\exp(L\delta)] \leq \sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{L^2 M_2}{2}\right)^n = \exp\left(\frac{L^2 \sigma^2}{2}\right)$$

and therefore

$$P(\delta \geq \lambda\sigma) \leq \exp\left(\frac{L^2 \sigma^2}{2} - L\lambda\sigma\right) \text{ for all positive } L.$$

The Chernoff Bound is now obtained by choosing  $L$  such that the right hand side of this expression is optimal. By differentiating it is easy to see that  $L = \lambda/\sigma$  is the correct choice.

It now follows immediately that

$$Pr(\delta \geq \lambda\sigma) \leq e^{-\lambda^2/2} \text{ and}$$

$$Pr(|\delta| = \lambda\sigma) \leq 2e^{-\lambda^2/2} \text{ as stated.}$$

The Chernoff Bound on the probability is much tighter than the one obtained by Chebyshev's Inequality, and is almost as good as having a normal distribution. For comparison we tabulate the  $\lambda$ 's obtained from each method:

$\Pr( \delta  \geq \lambda\sigma)$	$\lambda$ ( $\delta$ is normal)	Chernoff Bound for $\lambda$ ( $\delta$ is roundoff error)	Chebyshev Bound for $\lambda$ ( $\delta$ arbitrary)
.10	1.645	2.448	3.162
.05	1.960	2.716	4.472
.02	2.326	3.035	7.072
.01	2.576	3.255	10.00
.002798	3	3.625	18.90
.002	3.090	3.717	22.36
.001	3.290	3.899	31.62
.0002	3.719	4.292	70.71
.00002	4.265	4.799	223.6
.000002	4.753	5.257	707.2

Table 1. Bounds for  $\lambda$  given information about  $\delta$

Table 1 shows us that 99.7% confidence intervals on the magnitude of  $\delta$  are  $(-3.625\sigma, 3.625\sigma)$ ; this result compares quite closely with the  $(-3\sigma, +3\sigma)$  bounds we could use if it were known that  $\delta$  were normal, and has been used in the program that actually does statistical error analysis (Section 6).

## 5. CORRELATED ERRORS

Until now we have been assuming that all our relative errors were independent random variables. Now we consider the case in which they are not all independent, i.e., some of the errors are related to others in a common expression. This is the case, for example, in the evaluation of  $f_1(A \cdot A)$ . Assuming that

$$f_1(A) = A(1 + \delta_A) = [A, 1 * p_A] \quad (5.1)$$

then our prior analysis would claim

$$\begin{aligned} f_1(A \cdot A) &= A \cdot A(1 + \delta_A + \delta_A + \delta_R) \\ &= [A \cdot A, 1 * p_A * p_A * p_R]. \end{aligned} \quad (5.2)$$

This is not correct, however; it is too optimistic. The correct expression is

$$\begin{aligned} f_1(A \cdot A) &= A \cdot A(1 + 2\delta_A + \delta_R) \\ &= [A \cdot A, \Delta * p_A^{(2)} * p_R] \end{aligned} \quad (5.3)$$

where  $p_A^{(2)}$  is  $p_A$  stretched by the factor 2, as in (3.11); i.e.,

$$p_A^{(2)}(t) = \frac{1}{2} p_A(t/2). \quad (5.4)$$

We remark that the rounding error  $\delta_R$  (with density  $p_R$ ) is always initially an independent random variable. Correlation comes from repeated use of an existing error as in the case with  $\delta_A$  above; new rounding errors are always uncorrelated with anything.

The generalization of the above analysis to the general case is straightforward. Consider a smeared number  $\tilde{A}$

$$\tilde{A} = [A, 1 * p_1^{(c_1)} * p_2^{(c_2)} * \dots * p_N^{(c_N)}] \quad (5.5)$$

where  $p_k$  is a probability density ( $k = 1, \dots, N$ ) and the  $c_k$  are stretching factors. Suppose that two distinct error densities  $p_i$ ,  $p_j$  correspond to the same error variable  $\delta$ . Without loss of generality (convolution is commutative) we can assume these densities are  $p_1$  and  $p_2$ ;  $\tilde{A}$  can then be rewritten as

$$\tilde{A} = [A, 1 * p_1^{(c_1 + c_2)} * p_3^{(c_3)} * \dots * p_N^{(c_N)}] \quad (5.6)$$

This corresponds to merging the two references to  $\delta$  into one. Other references to  $\delta$  should also be merged.

Note that, because convolution of densities (or addition of random variables) is commutative, correlated errors must be merged throughout computation. For example, in computing

$$X = f1(\dots(((A + B) + C) + D) + E) + \dots + Z) + A) \quad (5.7)$$

the error  $\delta_A$  involved in the last step of the addition will be slightly correlated with the error of the rest of the sum; the resulting smeared number is

$$X = [X, 1 * p_A^{(2A/X)} * p_B^{(B/X)} * \dots * p_Z^{(Z/X)}] \quad (5.8)$$

where  $X = 2A + B + C + \dots + Z$ .

## 6. THE SNORT SYSTEM AND PROBLEMS OF IMPLEMENTATION

As this project of statistical error analysis progressed, it was hoped that the final results could be automated to perform a posteriori error bound calculation in the spirit of Interval Analysis. Since for small relative errors  $\delta = (f1(x) - x)/x$  we have

$$x \approx f1(x) (1 - \delta)$$

it seems natural to evaluate  $f1(x)$  in the natural way and use statistical bounds on  $\delta$  to bound the error in  $f1(x)$ . While Interval Analysis will give results like

$$x \in [x_{low}, x_{high}]$$

we would give results like

$$x \in [f1(x) (1 + \mu - B), f1(x) (1 + \mu + B)]$$

where  $B = (3.625\sigma)$  is the 99.7%-confidence bound on  $\delta$ ,  $\mu$  and  $\sigma^2$  being the computed mean and variance of  $\delta$  viewed as a random variable.

Automation of the statistical analysis was in fact completed by development of SNORTRAN, a PL/I-like language. It is similar in philosophy to PL/I-FORMAC in that SNORTRAN source is preprocessed and expanded into PL/I, which is then handled by the PL/I compiler. The preprocessor, a SNOBOL routine included at the end of this section, is unsophisticated: current need did not warrant sophistication.

The results of several test programs (also included at the end of this section) reflect negatively on continued use of statistical

relative error analysis for almost all problems. There are at least four reasons for the discouragement:

1. There is no way to represent 0 as a smeared number.
2. In subtraction of two numbers with similar values, cancellation causes enormous increase in error variance (stretching factors get large). One or two such cancellations will result in a value of  $\sigma$  so large that "useless" error bounds result.
3. In only the most well-conditioned, stable computations (e.g., repeated multiplication) are statistical bounds competitive with Interval Analysis.
4. Correlation between errors cannot be easily taken into account (as shown in Section 5) so statistical "bounds" ignoring correlation may use smaller stretching factors than they should, and give incorrectly small results.

The actual implementation uses Interval Analysis methods to generate worst-case values of  $\sigma$  and  $\mu$ , hence statistical bounds generated by the program should be bounds (ignoring correlation effects). There are four basic SNORT commands, all of them keyword-driven for simplicity:

```
INITIALIZE { ROUNDED }
           { CHOPPED };
SMEAR  <variable> [, <variable>];
EVAL   <variable> = <expression>;
PRINT  <variable> [, <variable>];
```



INITIALIZE causes the preprocessor to copy SNORT system declarations and routines into the translated program output, depending on the type of arithmetic being simulated. If rounded arithmetic is being used, the preprocessor copies in the file SNORT.RND, otherwise copying SNORT.CHP. The former is reproduced at the end of this section; the latter, requiring full Interval Analysis to compute error means, was never written.

SMEAR <variable> results in declarations for <variable> which permit it to be used as a smeared number. "<variable>" may be subscripted if desired (e.g., SMEAR A (10,0:5)) and a list of <variable>s generates a list of appropriate declarations.

EVAL <variable> = <expression> causes an SLR(1) parser to break down <expression> into SNORT system subroutine calls. Currently <expression> should only involve the operators +,-,/,\* although this could easily be expanded or modified by changing the parser tables. All variables in <variable> and <expression> should be SMEARed--the preprocessor does not currently check.

PRINT provides formatted output of smeared variables. It lists computed values,  $3.625\sigma$  bounds and the value bounds they imply, and Interval Analysis value bounds.

Four sample runs are included in the following pages, involving algorithms which highlight progressively more serious problems with statistical relative error analysis. The first shows that even in the best case of serial multiplication, statistical analysis does not give better bounds than interval analysis. The other three show



respectively what happens when addition, subtraction, and computed zeroes are added.

Preliminary example: Computation of  $\sin(1.0)$

The first test of the SNORTRAN system attempted to compute  $\sin(1) = .84147098\dots$  via the little-used relationship

$$\frac{\sin \theta}{\theta} = \prod_{n=1}^{\infty} \cos\left(\frac{\theta}{2^n}\right).$$

This concentrates on the roundoff properties of multiplication (ignoring computation of the cosines) which we know statistical analysis handles well. Unfortunately the results are not encouraging. Statistical analysis obtains the bounds

$$\sin(1) \in (.84146956, .84147235)$$

after 32 factors, while Interval Analysis gets

$$\sin(1) \in (.8414694, .8414710).$$

```

00100 PRELIM: PPOCEDURE OPTIONS(MAIN) REORDER;
00200 INITIALIZE ROUNDED;
00300 /*-----*/
00400 /* EVALUATE SIN(1.0) USING THE RELATIONSHIP */
00500 /* SIN(X)/X = COS(X/2) COS(X/4) COS(X/8) COS(X/16)....*/
00600 /* -- COSINES TREATED AS HAVING REPRESENTATION ERROR. */
00700 /* -- SERIES IS TRUNCATED AFTER 32 COSINE FACTORS. */
00800 /* -- CORRECT ANSWER IS 0.84147 09848 07896 50665... */
00900 /*-----*/
01000 SMEAR SIN_OF_A_BIT, COS_FACTOR;
01100 EVAL SIN_OF_A_BIT = 1.0EQ;
01200 DO N=1 TO 32;
01300     CALL SETCONSTANT( COS_FACTOR,
01400                       COS( 2.0000000000000000E0*(-N)) );
01500     EVAL SIN_OF_A_BIT = SIN_OF_A_BIT * COS_FACTOR;
01600     END;
01700 PRINT SIN_OF_A_BIT;
01800 END PPELIM;

SIN_OF_A_BIT      COMPUTED VALUE:      8.414709568023681E-01
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -2.779955685138702E+01 , 2.779955685138702E+01 ) * BETA**(-T)
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 8.414695624995408E-01 , 8.414723511051954E-01 )
INTERVAL ANALYSIS BOUNDS: ( 8.414694E-01 , 8.414710E-01 )

```

Preliminary Example SNORTRAN Source and Output

Example 0: Computation of  $\exp(8.0)$

This program involves the use of addition, multiplication, and division by computing  $\exp(8)$  with the usual power series, truncated at 40 terms. To ten digits  $e^8 = 2980.957987$ , and at the end statistical analysis claims

$$e^8 \in (2980.9549\dots, 2980.961\dots)$$

while Interval Analysis gets

$$e^8 \in (2980.955, 2980.958)$$

- a slightly tighter bound.

```

00100      EXO: PROCEDURE OPTIONS(MAIN) REORDER;
00200      INITIALIZE ROUNDED;
00300      SMEAR X, XN, FACTN, SN, EXPX;
00400      EVAL X = 8EO; /* EVALUATE EXP(8) THE HARD WAY */
00500      EVAL XN = 1.OEO;
00600      EVAL FACTN = 1.OEO;
00700      EVAL EXPX = 1.OEO;
00800      DO N=1 TO 40;
00900      CALL SETCONSTANT(SN,N);
01000      EVAL XN = XN*X;
01100      EVAL FACTN = FACTN * SN;
01200      EVAL EXPX = EXPX + XN/FACTN;
01300      PRINT EXPX;
01400      END;
01500      END EXO;

EXPX      COMPUTED VALUE: 2.980958007812500E+03
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -1.649799942970276E+01 , 1.649799942970276F+01 ) * BETA**(-1)
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 2.980955076465397E+03 , 2.980960939159603E+03 )
INTERVAL ANALYSIS BOUNDS: ( 2.980955E+03 , 2.980958E+03 )

EXPX      COMPUTED VALUE: 2.980958007812500E+03
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -1.678534758090973E+01 , 1.678534758090973E+01 ) * BETA**(-1)
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 2.980955025409681E+03 , 2.980960990215318E+03 )
INTERVAL ANALYSIS BOUNDS: ( 2.980955E+03 , 2.980958E+03 )

EXPX      COMPUTED VALUE: 2.980958007812500E+03
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -1.706785929203033E+01 , 1.706785929203033E+01 ) * BETA**(-1)
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 2.980954975213299E+03 , 2.980961040411701E+03 )
INTERVAL ANALYSIS BOUNDS: ( 2.980955E+03 , 2.980958E+03 )

EXPX      COMPUTED VALUE: 2.980958007812500E+03
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -1.734576964378357E+01 , 1.734576964378357E+01 ) * BETA**(-1)
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 2.980954925834482E+03 , 2.980961089790518E+03 )
INTERVAL ANALYSIS BOUNDS: ( 2.980955E+03 , 2.980958E+03 )

```

### Example 1: Computation of $\cos(4\pi)$

This problem involves all four arithmetic operations (using again the usual power series) and is not well-disposed to statistical analysis since there is eventually tremendous cancellation between terms. Note that to the first ten terms the computed sum is roughly -309, and the eleventh term is around +387.6, resulting in a factor of 16 increase in variance. Statistical analysis can no longer give decent results: it provides the bounds

$$\cos(4\pi) \in (-.4183606\dots, 2.4383515\dots)$$

whereas Interval Analysis gets

$$\cos(4\pi) \in (.7754054, 1.227057).$$

Neither of these results are good but the statistical bound is terrible. The relative error bound is so big that the assumption that second order effects are negligible is in question.

```

00100 EX1: PROCEDURE OPTICS(MAIN) REORDER;
00200 INITIALIZE ROUNDED;
00300 SMEAR X, X2, XN, FACTN, SN, COSX;
00400 EVAL X = 12.566370616E0; /* = 4*PI; WE ARE GOING TO */
00500 EVAL X2 = X*X; /* COMPUTE CUS(4*PI) THE HARD, */
00600 EVAL XN = 1.0E0; /* SLOWLY CONVERGING WAY. */
00700 EVAL FACTN = 1.0E0;
00800 EVAL COSX = 1.0E0;
00900 DO N=2 TO 40 BY 2;
01000 CALL SETCONSTANT(SN,N);
01100 EVAL XN = -XN*X2;
01200 EVAL FACTN = FACTN * SN * (SN-1.E0);
01300 EVAL COSX = COSX + XN/FACTN;
01400 PRINT COSX;
01500 END;
01600 END EX1;

COSX COMPUTED VALUE: -7.79568481453125E+01
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -7.925490617752075E+00 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( -7.795681131802388E+01 ;
INTERVAL ANALYSIS BOUNDS: ( -7.795682E+01 ;
7.925490617752075E+00 ) * BETA**(-T)
-7.795681131802388E+01 )
-7.795682E+01 )

COSX COMPUTED VALUE: 9.610737304687500E+02
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -1.100492358207703E+01 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 9.610743608798360E+02 ;
INTERVAL ANALYSIS BOUNDS: ( 9.610728E+02 ;
9.610737E+02 )
1.100492358207703E+01 ) * BETA**(-T)
9.610743608798360E+02 )
9.610737E+02 )

COSX COMPUTED VALUE: -4.5081646250000E+03
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -1.477807629108429E+01 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( -4.508168033480193E+03 ;
INTERVAL ANALYSIS BOUNDS: ( -4.508164E+03 ;
1.477807629108429E+01 ) * BETA**(-T)
-4.508168033480193E+03 )
-4.508164E+03 )

COSX COMPUTED VALUE: 1.091447265625000E+04
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -2.054601681232452E+01 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 1.091445928997167E+04 ;
INTERVAL ANALYSIS BOUNDS: ( 1.091445E+04 ;
1.091448E+04 )
2.054601681232452E+01 ) * BETA**(-T)
1.091445928997167E+04 )
1.091448E+04 )

COSX COMPUTED VALUE: -1.614603515625000E+04
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -2.946258699393951E+01 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( -1.61460351041588E+04 ;
INTERVAL ANALYSIS BOUNDS: ( -1.614606E+04 ;
1.614606E+04 )
2.946258699393951E+01 ) * BETA**(-T)
-1.61460351041588E+04 )
-1.614606E+04 )

COSX COMPUTED VALUE: 1.622687890625000E+04
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -4.505303657054901E+01 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 1.62268353107228E+04 ;
INTERVAL ANALYSIS BOUNDS: ( 1.622678E+04 ;
1.622694E+04 )
4.505303657054901E+01 ) * BETA**(-T)
1.62268353107228E+04 )
1.622694E+04 )

COSX COMPUTED VALUE: -1.186170703125000E+04
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -7.559576541259765E+01 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( -1.186176047658189E+04 ;
INTERVAL ANALYSIS BOUNDS: ( -1.186185E+04 ;
1.186185E+04 )
7.559576541259765E+01 ) * BETA**(-T)
-1.186176047658189E+04 )
-1.186185E+04 )

COSX COMPUTED VALUE: 6.619647656250000E+03
3.625*SIGMA BOUNDS ON RELATIVE ERROR: ( -1.461725826263428E+02 ;
IMPLIED 99.7%-CONFIDENCE BOUNDS: ( 6.619789980397908E+03 ;
INTERVAL ANALYSIS BOUNDS: ( 6.619652E+03 ;
6.620012E+03 )
1.461725826263428E+02 ) * BETA**(-T)
6.619789980397908E+03 )
6.620012E+03 )

```



```

COSX      COMPUTED VALUE:      -2.917707031250000E+03      ,      3.383616244334082E+02      ,      BETA**(-T)
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -3.383616244334082E+02      ,
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      -2.917707031250000E+03      ,
INTERVAL ANALYSIS BOUNDS:      (      -2.917707031250000E+03      ,

COSX      COMPUTED VALUE:      1.045739746093750E+03
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -9.47942240535156E+02      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      1.045739746093750E+03      ,
INTERVAL ANALYSIS BOUNDS:      (      1.045739746093750E+03      ,

COSX      COMPUTED VALUE:      -3.08983867187500E+02
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -3.21256068945313E+03      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      -3.08983867187500E+02      ,
INTERVAL ANALYSIS BOUNDS:      (      -3.08983867187500E+02      ,

COSX      COMPUTED VALUE:      7.850943593750000E+01
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -1.268161645507813E+04      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      7.850943593750000E+01      ,
INTERVAL ANALYSIS BOUNDS:      (      7.850943593750000E+01      ,

COSX      COMPUTED VALUE:      -1.558444213667188E+01
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -6.501621582031250E+04      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      -1.558444213667188E+01      ,
INTERVAL ANALYSIS BOUNDS:      (      -1.558444213667188E+01      ,

COSX      COMPUTED VALUE:      4.082443237304088E+00
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -2.672873359375000E+05      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      4.082443237304088E+00      ,
INTERVAL ANALYSIS BOUNDS:      (      4.082443237304088E+00      ,

COSX      COMPUTED VALUE:      5.127058029174804E-01
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -4.132032375000000E+06      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      5.127058029174804E-01      ,
INTERVAL ANALYSIS BOUNDS:      (      5.127058029174804E-01      ,

COSX      COMPUTED VALUE:      1.030962181091309E+00
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -3.562749687500000E+06      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      1.030962181091309E+00      ,
INTERVAL ANALYSIS BOUNDS:      (      1.030962181091309E+00      ,

COSX      COMPUTED VALUE:      1.000984191894531E+00
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -6.034146000000000E+06      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      1.000984191894531E+00      ,
INTERVAL ANALYSIS BOUNDS:      (      1.000984191894531E+00      ,

COSX      COMPUTED VALUE:      1.011007308959961E+00
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -9.466379375000000E+06      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      1.011007308959961E+00      ,
INTERVAL ANALYSIS BOUNDS:      (      1.011007308959961E+00      ,

COSX      COMPUTED VALUE:      1.009881973266602E+00
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -1.499482700000000E+07      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      1.009881973266602E+00      ,
INTERVAL ANALYSIS BOUNDS:      (      1.009881973266602E+00      ,

COSX      COMPUTED VALUE:      1.009995460510254E+00
3.625*SIGMA BOUNDS ON RELATIVE ERROR:      (      -2.372667987500000E+07      ,      BETA**(-T)
IMPLIED 99.74-CONFIDENCE BOUNDS:      (      1.009995460510254E+00      ,
INTERVAL ANALYSIS BOUNDS:      (      1.009995460510254E+00      ,

```

Example 1:  $\cos(4\pi)$  SNORTRAN Source and Output

### Final example: Lanczos' Algorithm

Lanczos' Algorithm is a method for reducing a symmetric matrix to tridiagonal form. Without re-orthogonalization, the vectors of the orthogonal transformation which it generates deviate gradually from orthonormality. The SNORTRAN coding of the algorithm was written before the problem with computed zeroes in statistical relative error analysis was fully appreciated, otherwise it would never have been done. It is included here as an example of how the SNORTRAN system was intended to be used (the PL/I preprocessor-output is included, since the algorithm uses so many SNORTRAN features), and also as an indicator of how severe a problem the irrepresentability of zero is. The program is initially given  $v_1 = (1, 0, 0, 0, \dots)$  and is asked to generate 63 other mutually orthonormal vectors. Statistical bounds for  $\langle v_1, v_2 \rangle$  are huge because the result is so close to zero, and during computation of  $\langle v_2, v_3 \rangle$  the program blows up because a zero is actually encountered.



Computation of  
the square root  
of the smeared  
number GAMMA

```

06100          CALL INNER(IPVAL(I),TEMP1,TEMP2);
06200          CALL SMEARED_PRINT(STRING(I),IPVAL(I));
06300          END;
06400      END;
06500
06600  AMULT:  PROC(U,V);
06700          SMEAR U(*), V(*);
06800          DECLARE (I,J,K) FIXED BIN;
06900          K = 0;
07000          DO J=1 TO P;
07100          DO I=1 TO N;
07200              K = K + 1;
07300              EVAL U(K) = 4 * V(K);
07400              IF I>1 THEN EVAL U(K) = U(K) - V(K-1);
07500              IF I<N THEN EVAL U(K) = U(K) - V(K+1);
07600              IF J>1 THEN EVAL U(K) = U(K) - V(K-N);
07700              IF J<P THEN EVAL U(K) = U(K) - V(K+N);
07800          END;
07900          END;
08000      END AMULT;
08100
08200
08300  INNER:  PROC(IP,U,V);
08400          SMEAR IP, U(*), V(*);
08500          DECLARE K FIXED BIN;
08600          IP = 0;
08700          DO K=1 TO NP;
08800              EVAL IP = IP + U(K) * V(K);
08900          END;
09000      END INNER;
09100
09200  END LANCZOS;

```

Final Example: Lanczos' Algorithm SNORTRAN Source

## Preprocessor Output of LANCZOS Program

MT LEV NT

```

1      0 | LANCZOS:  PROC OPTIONS(MAIN) REORDER;
2      1 0 |         DECLARE
          |             (I,J,K,II,IK,IX,IL,NLIMIT,N,P,NP) FIXED BIN,
          |             T  FLOAT BINARY(53);
3      1 0 | DECLARE
          |     1 U(64) LIKE SMEARED_NUMBER,
          |     1 V(8,64) LIKE SMEARED_NUMBER,
          |     1 W(64) LIKE SMEARED_NUMBER,
          |     1 ALPHA LIKE SMEARED_NUMBER,
          |     1 BETA LIKE SMEARED_NUMBER,
          |     1 GAMMA LIKE SMEARED_NUMBER;
4      1 0 | DECLARE
          |     1 VSAVE(8,64) LIKE SMEARED_NUMBER,
          |     1 IPVAL(8) LIKE SMEARED_NUMBER,
          |     1 TEMP1(64) LIKE SMEARED_NUMBER,
          |     1 TEMP2(64) LIKE SMEARED_NUMBER;
5      1 0 | DECLARE
          |             STRING (8) CHAR(13) INIT('<V(J),V(J)>','<V(J),V(J-1)>',
          |             '<V(J),V(J-2)>','<V(J),V(J-3)>','<V(J),V(J-4)>',
          |             '<V(J),V(J-5)>','<V(J),V(J-6)>','<V(J),V(J-7)>');
6      1 0 | DECLARE
          |     1 SMEARED_NUMBER,
          |     2 COMPUTED_VALUE  FLOAT BINARY(21),
          |     2 VARIANCE        FLOAT BINARY(21), /*UNITS BETA**(-2T)*/
          |     2 INTDATA,
          |     3 IA_LB           FLOAT BINARY(21), /* INTERVAL */
          |     3 IA_UB           FLOAT BINARY(21); /* ANALYSIS DATA */
7      1 0 | DECLARE
          |     1 REG(10) LIKE SMEARED_NUMBER,
          |     BETA_MT  FLOAT BINARY(21) INIT(1.0E-24B), /*16**(-6)*/
          |     HUGE     FLOAT BINARY(21) INIT(1.0E+88B);
8      1 0 | SMEARED_ADD:  PROC (A,B,C);
          |     /* PERFORMS ADDITION  A = B + C  FOR SMEARED NUMBERS  */
9      2 0 |     DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,
          |             (T,T1,T2,T3) FLOAT BINARY(53);
10     2 0 |     T = B.COMPUTED_VALUE;
11     2 0 |     A.COMPUTED_VALUE = ROUND(T + C.COMPUTED_VALUE);
12     2 0 |     T2 = MAXSQ(B.INTDATA);
13     2 0 |     T3 = MAXSQ(C.INTDATA);
14     2 0 |     CALL INTERVAL(A.INTDATA,1 /*ADD*/, B.INTDATA,C.INTDATA);
15     2 0 |     T1 = MINSQ(A.INTDATA);
16     2 0 |     IF ((T2+T3)/T1) > HUGE THEN
17     2 1 |         DO:  A.VARIANCE = HUGE; RETURN;  END;
20     2 0 |     A.VARIANCE = ROUNDUP( T2/T1*B.VARIANCE + T3/T1*C.VARIANCE );
21     2 0 |     CALL ADD_ROUNDING_ERROR(A);
22     2 0 |     END SMEARED_ADD;

```

```

23  1  0  | SMEARED_NEGATE: PROC (A,B);
          | /* PERFORMS NEGATION A = -B FOR SMEARED NUMBERS */
24  2  0  | DECLARE 1 (A,B) LIKE SMEARED_NUMBER,
          | T FLOAT BINARY(21);
25  2  0  | A.COMPUTED_VALUE = -B.COMPUTED_VALUE;
26  2  0  | A.VARIANCE = B.VARIANCE;
27  2  0  | T = B.IA_LB;
28  2  0  | A.IA_LB = -B.IA_UB;
29  2  0  | A.IA_UB = -T;
30  2  0  | END SMEARED_NEGATE;

31  1  0  | SMEARED_SUB: PROC (A,B,C);
          | /* PERFORMS SUBTRACTION A = B - C FOR SMEARED NUMBERS */
32  2  0  | DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,
          | (T,T1,T2,T3) FLOAT BINARY(53);
33  2  0  | T = B.COMPUTED_VALUE;
34  2  0  | A.COMPUTED_VALUE = ROUND(T - C.COMPUTED_VALUE);
35  2  0  | T2 = MAXSQ(B.INTDATA);
36  2  0  | T3 = MAXSQ(C.INTDATA);
37  2  0  | CALL INTERVAL(A.INTDATA,2 /*SUB*/, B.INTDATA,C.INTDATA);
38  2  0  | T1 = MINSQ(A.INTDATA);
39  2  0  | IF ((T2+T3)/T1) > HUGE THEN
40  2  1  | DO; A.VARIANCE = HUGE; RETURN; END;
43  2  0  | A.VARIANCE = ROUNDUP( T2/T1*B.VARIANCE + T3/T1*C.VARIANCE );
44  2  0  | CALL ADD_ROUNDING_ERROR(A);
45  2  0  | END SMEARED_SUB;

46  1  0  | SMEARED_MULT: PROC(A,B,C);
          | /* PERFORMS MULTIPLICATION A = B * C FOR SMEARED NUMBERS */
47  2  0  | DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,
          | T FLOAT BINARY(53);
48  2  0  | T = B.COMPUTED_VALUE;
49  2  0  | A.COMPUTED_VALUE = ROUND(T * C.COMPUTED_VALUE);
50  2  0  | CALL INTERVAL(A.INTDATA,3 /*MULT*/,B.INTDATA,C.INTDATA);
51  2  0  | T = B.VARIANCE;
52  2  0  | A.VARIANCE = ROUNDUP(T + C.VARIANCE);
53  2  0  | CALL ADD_ROUNDING_ERROR(A);
54  2  0  | END SMEARED_MULT;

55  1  0  | SMEARED_DIV: PROC(A,B,C);
          | /* PERFORMS DIVISION A = B / C FOR SMEARED NUMBERS */
56  2  0  | DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,
          | T FLOAT BINARY(53);
57  2  0  | T = B.COMPUTED_VALUE;
58  2  0  | A.COMPUTED_VALUE = ROUND(T / C.COMPUTED_VALUE);
59  2  0  | CALL INTERVAL(A.INTDATA,4 /*DIV*/,B.INTDATA,C.INTDATA);
60  2  0  | T = B.VARIANCE;
61  2  0  | A.VARIANCE = ROUNDUP(T + C.VARIANCE);
62  2  0  | CALL ADD_ROUNDING_ERROR(A);
63  2  0  | END SMEARED_DIV;

64  1  0  | SETCONSTANT: PROC(A,C);
          | /* SETS SMEARED NUMBER A TO THE INITIAL VALUE C */
65  2  0  | DECLARE 1 A LIKE SMEARED_NUMBER,
          | C FLOAT BINARY(53);
66  2  0  | A.COMPUTED_VALUE = ROUND(C);
67  2  0  | A.IA_LB = ROUNDDOWN(C);
68  2  0  | A.IA_UB = ROUNDUP(C);
69  2  0  | A.VARIANCE = OE0B;
70  2  0  | IF C /= PRECISION(C,21) THEN
          | CALL ADD_ROUNDING_ERROR(A);
71  2  0  | END SETCONSTANT;

```

```

72 1 0 | INTERVAL: PROC(A,OP,B,C);
      /* PERFORMS INTERVAL ARITHMETIC   A = B OP C   */
73 2 0 | DECLARE 1 INTERVAL,
      2 IA_LB      FLOAT BINARY(21),
      2 IA_UB      FLOAT BINARY(21),
      1 (A,B,C) LIKE INTERVAL,
      (S1,S2,S3,S4) FLOAT BINARY(53),
      IOP(4) LABEL,
      OP          FIXED BINARY(15);
74 2 0 | S1 = B.IA_LB; S2 = B.IA_UB; S3 = C.IA_LB; S4 = C.IA_UB;
78 2 0 | GO TO IOP(OP);

79 2 0 | IOP(1):      /* ADDITION */
      A.IA_LB = ROUNDDOWN(S1 + S3);
80 2 0 | A.IA_UB = ROUNDUP (S2 + S4);
81 2 0 | RETURN;

82 2 0 | IOP(2):      /* SUBTRACTION */
      A.IA_LB = ROUNDDOWN(S1 - S4);
83 2 0 | A.IA_UB = ROUNDUP (S2 - S3);
84 2 0 | RETURN;

85 2 0 | IOP(3):      /* MULTIPLICATION */
      A.IA_LB = MIN( ROUNDDOWN(S1 * S3),
                    ROUNDDOWN(S1 * S4),
                    ROUNDDOWN(S2 * S3),
                    ROUNDDOWN(S2 * S4) );
86 2 0 | A.IA_UB = MAX( ROUNDUP (S1 * S3),
                    ROUNDUP (S1 * S4),
                    ROUNDUP (S2 * S3),
                    ROUNDUP (S2 * S4) );

87 2 0 | RETURN;

88 2 0 | IOP(4):      /* DIVISION */
      A.IA_LB = MIN( ROUNDDOWN(S1 / S3),
                    ROUNDDOWN(S1 / S4),
                    ROUNDDOWN(S2 / S3),
                    ROUNDDOWN(S2 / S4) );
89 2 0 | A.IA_UB = MAX( ROUNDUP (S1 / S3),
                    ROUNDUP (S1 / S4),
                    ROUNDUP (S2 / S3),
                    ROUNDUP (S2 / S4) );

90 2 0 | RETURN;

91 2 0 | END INTERVAL;

92 1 0 | MAXSQ: PROC(A) RETURNS(FLOAT BINARY(53));
      /* RETURNS MAX VALUE OF A**2 OVER INTERVAL A */
93 2 0 | DECLARE 1 A, /* INTERVAL */
      2 LB FLOAT BINARY(21),
      2 UB FLOAT BINARY(21),
      (T1,T2) FLOAT BINARY(53);
94 2 0 | T1 = PRECISION(A.LB,53)**2;
95 2 0 | T2 = PRECISION(A.UB,53)**2;
96 2 0 | T1 = MAX(T1,T2);
97 2 0 | RETURN(T1);
98 2 0 | END MAXSQ;

```



```

99 1 0 | MINSQ: PROC(A) RETURNS(FLOAT BINARY(53));
      | /* RETURNS MIN VALUE OF A**2 OVER INTERVAL A */
100 2 0 | DECLARE 1 A, /* INTERVAL */
      | 2 LB FLOAT BINARY(21),
      | 2 UB FLOAT BINARY(21),
      | (T1,T2) FLOAT BINARY(53);
101 2 0 | IF (A.LB*A.UB <= OE0B) THEN /* INTERVAL CONTAINS ZERO */
102 2 1 | DO; T1 = 1EOB/HUGE; RETURN(T1); END;
105 2 0 | T1 = PRECISION(A.LB,53)**2;
106 2 0 | T2 = PRECISION(A.UB,53)**2;
107 2 0 | T1 = MIN(T1,T2);
108 2 0 | RETURN(T1);
109 2 0 | END MINSQ;

110 1 0 | ADD_ROUNDING_ERROR: PROC(A);
111 2 0 | DECLARE 1 A LIKE SMEARED_NUMBER,
      | (D,S) FLOAT BINARY(53);
      | /* COMPUTE D = MIN(MANTISSA(A)) */
      | /* S = 1 / (2D)**2 / 3 */
      | /* = COMPUTATION ROUNDOFF ERROR VARIANCE, */
      | /* UNITS BETA**(-2T). */
112 2 0 | IF A.VARIANCE >= HUGE | A.COMPUTED_VALUE=OE0B THEN RETURN;
113 2 0 | D = MIN(MANTISSA(A.IA_LB),MANTISSA(A.IA_UB));
114 2 0 | S = 1EOB / (D*D*12EO);
115 2 0 | A.VARIANCE = ROUNDUP(A.VARIANCE + D);
116 2 0 | END ADD_ROUNDING_ERROR;

117 1 0 | MANTISSA: PROC(X) RETURNS(FLOAT BINARY(53));
118 2 0 | DECLARE (X,Y) FLOAT BINARY(21);
119 2 0 | Y = ABS(X);
120 2 0 | DO WHILE (Y > 1.000000000000000000000000E0B);
121 2 1 | Y = Y * 1.000000000000000000000000E-4B;
122 2 1 | END;
123 2 0 | DO WHILE (Y < 1.000000000000000000000000E-4B);
124 2 1 | Y = Y * 1.000000000000000000000000E+4B;
125 2 1 | END;
126 2 0 | RETURN(Y);
127 2 0 | END MANTISSA;

128 1 0 | SMEARED_PRINT: PROC(AA,A);
129 2 0 | DECLARE 1 ALIKE SMEARED_NUMBER,
      | AA CHAR(*),
      | (S,S1,S2,ABS1,ABS2,ABS3,ABS4) FLOAT BIN(53);
130 2 0 | S = SQRT(A.VARIANCE);
131 2 0 | S1 = -3.625 * S; /* 99.7% CONFIDENCE LIMITS */
132 2 0 | S2 = +3.625 * S;
133 2 0 | ABS3 = A.COMPUTED_VALUE * (1 + S1*BETA_MT);
134 2 0 | ABS4 = A.COMPUTED_VALUE * (1 + S2*BETA_MT);
135 2 0 | IF A.COMPUTED_VALUE < 0 THEN
136 2 1 | DO; ABS1=ABS3; ABS3=ABS4; ABS4=ABS1; END;

140 2 0 | PUT FILE(SYSPRINT) SKIP EDIT
      | (AA,'COMPUTED VALUE:',A.COMPUTED_VALUE)
      | (SKIP,COL(7),A,X(8),A,E(25,15,16))
      | ('3.625*SIGMA BOUNDS ON RELATIVE ERROR:',(' ',S1,' ',S2,
      | ' ) * BETA**(-T)' )
      | (SKIP,COL(12),A,COL(50),A,2 (E(25,15,16),A))

```

```

/*                                MAIN LOOP                                */
DO IX=2 TO NP;
  IK = MOD(IX-2,8) + 1;
  DO II=1 TO NP;    TEMP1(II) = V(IK,II);    END;
  CALL INNER(ALPHA,TEMP1,U);
  DO K=1 TO NP;

```



```

187 1 2 | CALL SMEARED_MULT(REG(1),ALPHA,V(IK,K));
188 1 2 | CALL SMEARED_SUB(W(K),U(K),REG(1));
189 1 2 | END;
190 1 1 | CALL INNER(GAMMA,W,W);
191 1 1 | GAMMA.COMPUTED_VALUE=SQRT(GAMMA.COMPUTED_VALUE);
192 1 1 | GAMMA.VARIANCE = GAMMA.VARIANCE/2;
193 1 1 | T=GAMMA.IA_LB; GAMMA.IA_LB=ROUNDDOWN(SQRT(T));
195 1 1 | T=GAMMA.IA_UB; GAMMA.IA_UB=ROUNDUP (SQRT(T));
197 1 1 | IL = MOD(IK,8) + 1;
198 1 1 | DO K=1 TO NP;

199 1 2 | CALL SMEARED_DIV(V(IL,K),W(K),GAMMA);
200 1 2 | END;

201 1 1 | BETA = GAMMA;

202 1 1 | DO II=1 TO NP; TEMP1(II) = V(IL,II); END;
205 1 1 | CALL AMULT(U,TEMP1);
206 1 1 | DO K=1 TO NP;

207 1 2 | CALL SMEARED_MULT(REG(1),BETA,V(IK,K));
208 1 2 | CALL SMEARED_SUB(U(K),U(K),REG(1));
209 1 2 | END;
    /* END OF MAIN LOOP COMPUTATION */
    /* NOW CHECK INNER PRODUCTS... */
210 1 1 | IF IX<8 THEN NLIMIT=IX; ELSE NLIMIT=8;
212 1 1 | PUT SKIP EDIT ('J = ',IX) (SKIP(2),A,F(3));
213 1 1 | DO I=1 TO NLIMIT;
214 1 2 | J = MOD(IL+8-I,8) + 1;
215 1 2 | DO II=1 TO NP; TEMP1(II) = V(IL,II); END;
218 1 2 | DO II=1 TO NP; TEMP2(II) = V(J,II); END;
221 1 2 | CALL INNER(IPVAL(I),TEMP1,TEMP2);
222 1 2 | CALL SMEARED_PRINT(STRING(I),IPVAL(I));
223 1 2 | END;
224 1 1 | END;
225 1 0 | AMULT: PROC(U,V);

226 2 0 | DECLARE
    1 U(*) LIKE SMEARED_NUMBER,
    1 V(*) LIKE SMEARED_NUMBER;
227 2 0 | DECLARE (I,J,K) FIXED BIN;
228 2 0 | K = 0;
229 2 0 | DO J=1 TO P;
230 2 1 | DO I=1 TO N;
231 2 2 | K = K + 1;

232 2 2 | CALL SETCONSTANT(REG(1),4);
233 2 2 | CALL SMEARED_MULT(U(K),REG(1),V(K));
234 2 2 | IF I>1 THEN
    CALL SMEARED_SUB(U(K),U(K),V(K-1));
235 2 2 | IF I<N THEN
    CALL SMEARED_SUB(U(K),U(K),V(K+1));
236 2 2 | IF J>1 THEN
    CALL SMEARED_SUB(U(K),U(K),V(K-N));
237 2 2 | IF J<P THEN
    CALL SMEARED_SUB(U(K),U(K),V(K+N));
238 2 2 | END;
239 2 1 | END;
240 2 0 | END AMULT;

```



## SNORTRAN System Listings

1. SNORTRAN Preprocessor (SNOBOL) pp.49-55
2. SNORTRAN System Routines (PL/I) pp.56-60

(Referenced as file 'SNORT.RND' in line 09000 of preprocessor)

```

00100  *-----*
00200  *
00300  *           S N O R T R A N
00400  *
00500  *       A POSTERIORI ROUNDOFF ERROR LANGUAGE PREPROCESSOR
00600  *
00700  *-----*
00800      &TRACE = 1000
00900      DEFINE('FILL(VECTOR,STRING)I,VAL')
01000      DEFINE('PARSE.EXPR()I,J,LOOKAHEAD.TOKEN,LOOKAHEAD.SYMBOL')
01100      DEFINE('SEMANTICS(PCDNUM)')
01200      DEFINE('GETCHAR()')
01300      DEFINE('GETNONBLANK()')
01400      DEFINE('SCAN()')
01500      DEFINE('ADD()')
01600      DEFINE('REMOVE()')
01700      DEFINE('EMIT(STRING)')
01800      DEFINE('LASTEMIT(TARGET)X,Y')
01900      DEFINE('ALLOCATED.REG()I')
02000      DEFINE('FREE.REGISTER()I')
02100      DEFINE('CLEARREGS()I')
02200      KEYWORD = ('INITIALIZE' | 'SMEAR' | 'EVAL' | 'PRINT') . KEY ' '
02300      NREGS = 10
02400      REGISTER = ARRAY(NREGS)
02500      SEMICOLON = 1
02600      VARIABLE = 2
02700      EQUALS = 3
02800      PLUS = 4
02900      MINUS = 5
03000      STAR = 6
03100      SLASH = 7
03200      LPAREN = 8
03300      RPAREN = 9
03400      FACTOR = 10
03500      TERM = 11
03600      EXPR = 12
03700      ASSIGN = 13
03800      ACCEPT = 14
03900      COMMA = 20
04000  -EJECT
04100  *
04200  *       SLR(1) EXPRESSION PARSER TABLES
04300  *
04400      DEPTH = 50
04500      STATESTACK = ARRAY(DEPTH)
04600      SYMBOLSTACK = ARRAY(DEPTH)
04700      TOKENSTACK = ARRAY(DEPTH)
04800      STORED.DEPTH = 5
04900      STORED.TOKEN = ARRAY(STORED.DEPTH)
05000      STORED.SYMBOL = ARRAY(STORED.DEPTH)
05100      NPRJDS = 12; NTRANS = 76; NSTATES = 23
05200      TS = '2,13,3,1,2,4,5,8,10,11,12,2,8,10,11,2,,'
05300  +      '4,5,8,10,11,12,1,4,5,8,7,9,1,4,5,6,7,'
05400  +      '9,1,4,5,6,7,9,4,5,9,2,8,10,2,3,10,2,8,10,11,'
05500  +      '2,8,10,11,1,4,5,6,7,9,1,4,5,6,7,9,'
05600      TA = '2,3,4,0,5,6,7,8,9,10,11,5,8,9,12,5,8,9,13,5,'
05700  +      '6,7,8,9,10,14,-5,-5,-5,15,10,-5,-2,17,18,-6,-6,-6,15,16,'
05800  +      '-6,-7,-7,-7,15,10,-7,17,18,19,5,8,20,5,8,21,5,8,9,22,'
05900  +      '5,8,9,23,-3,-3,-3,15,16,-3,-4,-4,-4,15,16,-4,'
06000      FT = '1,3,4,5,0,12,16,20,0,27,33,36,42,48,51,54,57,61,0,0,'

```

```

06100 +      '0,65,71,'
06200 LT = '2,3,4,11,-12,15,19,26,-10,32,35,41,47,50,53,56,60,64,-11'
06300 +      ',-8,-9,70,76,'
06400 LS = '14,13,12,12,12,12,11,11,11,10,10,'
06500 RL = '2,3,3,3,1,2,2,3,3,1,3,1,'
06600 TRANS.SYMBOL = ARRAY(NTRANS); FILL(.TRANS.SYMBOL,TS)
06700 TRANS.ACTION = ARRAY(NTRANS); FILL(.TRANS.ACTION,TA)
06800 FIRST.TRANS = ARRAY(NSTATES); FILL(.FIRST.TRANS,FT)
06900 LAST.TRANS = ARRAY(NSTATES); FILL(.LAST.TRANS,LT)
07000 LHS.SYMBOL = ARRAY(NPRODS); FILL(.LHS.SYMBOL,LS)
07100 RHS.LENGTH = ARRAY(NPRODS); FILL(.RHS.LENGTH,RL)
07200 +                                     : (READ)
07300 *
07400 FILL      I = I + 1
07500 STRING SPAN('0123456789-') . VAL ',' = :F(RETURN)
07600 ITEM($VECTOR,I) = CONVERT(VAL, .INTEGER) : (FILL)
07700 *
07800 *
07900 READ      CARD = INPUT                                     :F(END)
08000 RESUME    CARD @P KEYWORD                                   :S(PREPARE)
08100          OUTPUT = DIFFER(TRIM(CARD)) CARD                  : (READ)
08200 *
08300 PREPARE   CARD LEN(P) . SNARK =
08400          OUTPUT = DIFFER(TRIM(SNARK)) ' ' SNARK : ($KEY)
08500 *
08600 * INITIALIZE -- COPY IN SYSTEM DECLARES AND SUBROUTINES
08700 *              DETERMINE ALSO ROUNDING MODE.
08800 *
08900 INITIALIZE SCAN()
09000          INITFILE = 'SNORT.RND'
09100          INITFILE = IDENT(SCAN(), 'CHOPPED') 'SNORT.CHP'
09200          INPUT(.INIT, INITFILE)
09300 INITCOPY  OUTPUT = INIT                                     :S(INITCOPY)
09400          SCAN()
09500          OUTPUT = NE(SCAN(TYPE, SEMICOLON) ' /***** ERROR IN INITIALIZE'
09600 +          ' STATEMENT *****/'                               : (RESUME)
09700 *
09800 * SMEAR -- GENERATE DECLARATIONS FOR SMEARED VARIABLES
09900 *
10000 SMEAR    SCAN()
10100          OUTPUT = ' DECLARE'
10200 SMEARLOOP  V = SCAN()
10300          EQ(SCANTYPE, VARIABLE)                               :F(SMERRCR)
10400          SCAN()
10500          EQ(SCANTYPE, COMMA)                                   :F(SMEARFIN)
10600          OUTPUT = '      1 ' V ' LIKE SMEARED_NUMBER;' : (SMEARLOOP)
10700 SMEARFIN  EQ(SCANTYPE, SEMICOLON)                             :F(SMERROR)
10800 SMEND     OUTPUT = '      1 ' V ' LIKE SMEARED_NUMBER;' : (RESUME)
10900 SMEARCR   OUTPUT = ' /***** ERRCR IN SMEAR LIST. END FORCED *****/'
11000 +                                     : (RESUME)
11100 *
11200 * PRINT -- FORMATTED PRINT OF SMEARED VARIABLES
11300 *
11400 PRINT     SCAN()
11500 PRLIST    V = SCAN()
11600          EQ(SCANTYPE, VARIABLE)                               :F(PRERROR)
11700          OUTPUT = " CALL SMEARED_PRINT(' ' V ' ',' ' V ');"
11800          SCAN()
11900          EQ(SCANTYPE, COMMA)                                   :S(PRLIST)
12000          EQ(SCANTYPE, SEMICOLON)                             :S(RESUME)

```



```

12100 PRERROR OUTPUT = ' /***** ERROR IN PRINT LIST.  END FORCED *****/'
12200                                     :(RESUME)
12300 -EJECT
12400 *
12500 *  EVAL      --  PARSE EXPRESSIONS OF SMEARED NUMBERS AND
12600 *              GENERATE APPROPRIATE SUBROUTINE CALLS.
12700 *
12800 EVAL      SCAN()
12900          PARSE.EXPR()                                     :(RESUME)
13000 *
13100 SEMANTICS                                             :(('$('PROD' PRODNUM))
13200 *-----*
13300 *  1:          S ==> <ASSIGNMENT> ;                                     *
13400 *-----*
13500 PRGD1                                             :(RETURN)
13600 *-----*
13700 *  2:          <ASSIGNMENT> ==> <VAR> = <EXPR>                                     *
13800 *-----*
13900 PRGD2  LASTEMIT(TOKENSTACK<STACKPTR>)
14000          CLEARREGS()                                     :(RETURN)
14100 *-----*
14200 *  3:          <EXPR> ==> <EXPR> + <TERM>                                     *
14300 *-----*
14400 PRGD3  FREE.REGISTER(TOKENSTACK<STACKPTR>)
14500          FREE.REGISTER(TOKENSTACK<STACKPTR + 2>)
14600          AREG = ALLOCATED.REG()
14700          EMIT( ' CALL SMEARED_ADD(' AREG ', ' TOKENSTACK<STACKPTR>
14800 +              ', ' TOKENSTACK<STACKPTR + 2> '); ' )
14900          TOKENSTACK<STACKPTR> = AREG                     :(RETURN)
15000 *-----*
15100 *  4:          <EXPR> ==> <EXPR> - <TERM>                                     *
15200 *-----*
15300 PRGD4  FREE.REGISTER(TOKENSTACK<STACKPTR>)
15400          FREE.REGISTER(TOKENSTACK<STACKPTR + 2>)
15500          AREG = ALLOCATED.REG()
15600          EMIT( ' CALL SMEARED_SUB(' AREG ', ' TOKENSTACK<STACKPTR>
15700 +              ', ' TOKENSTACK<STACKPTR + 2> '); ' )
15800          TOKENSTACK<STACKPTR> = AREG                     :(RETURN)
15900 *-----*
16000 *  5:          <EXPR> ==> <TERM>                                     *
16100 *-----*
16200 PRGD5                                             :(RETURN)
16300 *-----*
16400 *  6:          <EXPR> ==> + <TERM>                                     *
16500 *-----*
16600 PRGD6  TOKENSTACK<STACKPTR> = TOKENSTACK<STACKPTR + 1> :(RETURN)
16700 *-----*
16800 *  7:          <EXPR> ==> - <TERM>                                     *
16900 *-----*
17000 PRGD7  FREE.REGISTER(TOKENSTACK<STACKPTR + 1>)
17100          AREG = ALLOCATED.REG()
17200          EMIT( ' CALL SMEARED_NEGATE(' AREG ', '
17300 +              TOKENSTACK<STACKPTR + 1> '); ' )
17400          TOKENSTACK<STACKPTR> = AREG                     :(RETURN)
17500 *-----*
17600 *  8:          <TERM> ==> <TERM> * <FACTOR>                                     *
17700 *-----*
17800 PRGD8  FREE.REGISTER(TOKENSTACK<STACKPTR>)
17900          FREE.REGISTER(TOKENSTACK<STACKPTR + 2>)
18000          AREG = ALLOCATED.REG()

```

```

18100      EMIT( ' CALL SMEARED_MULT(' AREG ', ' TOKENSTACK<STACKPTR>
18200      +      ', ' TOKENSTACK<STACKPTR + 2> ');' )
18300      TOKENSTACK<STACKPTR> = AREG      : (RETURN)
18400      *-----*
18500      * 9:      <TERM> ==> <TERM> / <FACTOR>      *
18600      *-----*
18700      PROD9  FREE.REGISTER(TOKENSTACK<STACKPTR>)
18800      FREE.REGISTER(TOKENSTACK<STACKPTR + 2>)
18900      AREG = ALLOCATED.REG()
19000      EMIT( ' CALL SMEARED_DIV(' AREG ', ' TOKENSTACK<STACKPTR>
19100      +      ', ' TOKENSTACK<STACKPTR + 2> ');' )
19200      TOKENSTACK<STACKPTR> = AREG      : (RETURN)
19300      *-----*
19400      * 10:     <TERM> ==> <FACTOR>      *
19500      *-----*
19600      PROD10      : (RETURN)
19700      *-----*
19800      * 11:     <FACTOR> ==> { <EXPR> }      *
19900      *-----*
20000      PROD11  TOKENSTACK<STACKPTR> = TOKENSTACK<STACKPTR + 1> : (RETURN)
20100      *-----*
20200      *      <FACTOR> ==> <VAR>      *
20300      *-----*
20400      PROD12  EQ(VARINFO,0)      :S (RETURN)
20500      CCNSTANT.VAR  AREG = ALLOCATED.REG()
20600      EMIT( ' CALL SETCONSTANT(' AREG ', '
20700      +      TOKENSTACK<STACKPTR> ');' )
20800      TOKENSTACK<STACKPTR> = AREG      : (RETURN)
20900      *-----*
21000      ALLOCATED.REG  I = LT(I,NREGS) I + 1      :F (ALL.USED)
21100      REGISTER<I> = IDENT(REGISTER<I>) 'X'      :F (ALLOCATED.REG)
21200      ALLOCATE  ALLOCATED.REG = 'REG(' I ') '      : (RETURN)
21300      ALL.USED  OUTPUT = ' /***** ALL REGISTERS USED *****/' : (ALLOCATE)
21400      *
21500      FREE.REGISTER  REGNAME 'REG(' BREAK(')') . I :F (RETURN)
21600      REGISTER<I> =      : (RETURN)
21700      *
21800      CLEARREGS      I = LT(I,NREGS) I + 1      :F (RETURN)
21900      REGISTER<I> =      : (CLEARREGS)
22000      *
22100      EMIT  OUTPUT = DIFFER(SAVED.STRING) SAVED.STRING
22200      SAVED.STRING = STRING      : (RETURN)
22300      *
22400      LASTEMIT  SAVED.STRING (BREAK('(' LEN(1)) . X BREAK(')') =
22500      +      X      TARGET      :F (ASSMT)
22600      OUTPUT = SAVED.STRING
22700      SAVED.STRING =      : (RETURN)
22800      ASSMT  OUTPUT = ' ' TARGET ' = ' TOKENSTACK<STACKPTR + 2> '; '
22900      +      : (RETURN)
23000      -EJECT
23100      *-----*
23200      *
23300      *      SLR(1) EXPRESSION PARSER
23400      *
23500      *-----*
23600      PARSE.EXPR  STATESTACK<1> = 1
23700      SYMBOLSTACK<1> = 0
23800      CURRENT.STATE = 1
23900      STACKPTR = 2
24000      NEEDREAD = 'YES'

```



```

24100 *
24200 PARSE      I = FIRST.TRANS<CURRENT.STATE>
24300           J = LAST.TRANS<CURRENT.STATE>
24400           EQ(I,0)                                :S(REDUCE.LR0)
24500           NEEDREAD = IDENT(NEEDREAD) 'YES'      :S(FIND)
24600           LOOKAHEAD.TOKEN = SCAN()
24700           LOOKAHEAD.SYMBOL = SCANTYPE
24800 FIND      EQ(LOOKAHEAD.SYMBOL,TRANS.SYMBOL<I>)   :S(FOUND)
24900           I = LT(I,J) I + 1                      :S(FIND)F(SYNTAXERR)
25000 FOUND     J = TRANS.ACTION<I>
25100           EQ(J,0)                                :S(ACCEPT)
25200           LT(J,0)                                :S(REDUCE.LR1)
25300 *
25400 SHIFT     STATESTACK<STACKPTR> = J
25500           SYMBOLSTACK<STACKPTR> = LOOKAHEAD.SYMBOL
25600           TOKENSTACK<STACKPTR> = LOOKAHEAD.TOKEN
25700           CURRENT.STATE = J
25800           STACKPTR = LT(STACKPTR,DEPTH) STACKPTR + 1 :S(PARSE)
25900           OUTPUT = ' /***** STACK OVERFLOW *****/' :S(PARSE)
26000 *
26100 REDUCE.LR0  NEEDREAD = 'YES'                      :S(REDUCE)
26200 REDUCE.LR1  NEEDREAD =
26300 REDUCE      PROD = -J
26400           STACKPTR = STACKPTR - RHS.LENGTH<PROD>
26500           CURRENT.STATE = STATESTACK<STACKPTR - 1>
26600           SYMBOLSTACK<STACKPTR> = LHS.SYMBOL<PROD>
26700           I = FIRST.TRANS<CURRENT.STATE>
26800           J = LAST.TRANS<CURRENT.STATE>
26900 GOTO      EQ(SYMBOLSTACK<STACKPTR>,TRANS.SYMBOL<I>) :S(PUSHGOTO)
27000           I = LT(I,J) I + 1                      :S(GOTO)F(SYNTAXERR)
27100 PUSHGOTO   J = TRANS.ACTION<I>
27200           STATESTACK<STACKPTR> = J
27300           CURRENT.STATE = J
27400           SEMANTICS(PROD)
27500           STACKPTR = STACKPTR + 1                  :S(PARSE)
27600 *
27700 ACCEPT                                           :S(RETURN)
27800 SYNTAXERR OUTPUT = ' /***** SYNTAX ERROR *****/' :S(FRETURN)
27900 *
28000 *-----*
28100 -EJECT
28200 *
28300 *          SCAN -- LEXICAL ANALYZER.      RETURNED VALUE IS TOKEN STRING.
28400 *          ALSO RETURNS  SCANTYPE -- NUMERICAL TOKEN CODE
28500 *          VARINFO -- EXTRA INFO FOR <VAR> TOKENS
28600 *
28700 SCAN      SCAN =
28800           NEXTCHAR =
28900           GETNONBLANK()                            :S('CLASS' CLASSN))
29000 *-----*
29100 *          CLASS 1 -- LETTERS, BREAK CHARACTERS:  BUILD VARIABLES
29200 *-----*
29300 CLASS1     SCANTYPE = VARIABLE;  VARINFO = 0
29400 CLASS1A    ADD();  GETCHAR()
29500           GE(CLASSN,1) LE(CLASSN,2)                :S(CLASS1A)
29600 CLASS1B    IDENT(NEXTCHAR,' ')                    :F(CLASS1C)
29700           GETNONBLANK()
29800 CLASS1C    EQ(CLASSN,8)                            :F(OUT)
29900           PARENCOUNT = 1
30000 CLASS1D  ADD();  GETCHAR()

```

```

30100      PARENCOUNT = EQ(CLASSN,8) PARENCOUNT + 1 :S(CLASS1D)
30200      PARENCOUNT = EQ(CLASSN,9) PARENCOUNT - 1 :F(CLASS1D)
30300      EQ(PARENCOUNT,0)                               :F(CLASS1D)
30400      ADD()                                           :(OUT)
30500    *-----*
30600    *          CLASS 2 -- DIGITS: BUILD CONSTANTS
30700    *-----*
30800    CLASS2 SCANTYPE = VARIABLE; VARINFO = 1
30900    CLASS2A ADD(); GETNONBLANK()
31000           VARINFO = VARINFO + 1
31100           EQ(CLASSN,2)                                :S(CLASS2A)
31200           VARINFO = IDENT(NEXTCHAR,'.') 100            :S(CLASS2A)
31300           IDENT(NEXTCHAR,'E')                          :F(OUT)
31400           ADD(); GETCHAR()
31500           GE(CLASSN,2) LE(CLASSN,4)                    :F(OUT)
31600    CLASS2B ADD(); GETCHAR()
31700           EQ(CLASSN,2)                                    :S(CLASS2B)F(OUT)
31800    *-----*
31900    *          CLASSES 3-10 -- TERMINALS       3+ 4- 5/ 6* 7= 8( 9) 10;
32000    *-----*
32100    CLASS3 ADD(); GETCHAR(); SCANTYPE = PLUS
32200           DIFFER(UNARY) EQ(CLASSN,2)                     :S(CLASS2)F(OUT)
32300    CLASS4 ACD(); GETCHAR(); SCANTYPE = MINUS
32400           DIFFER(UNARY) EQ(CLASSN,2)                     :S(CLASS2)F(CUT)
32500    CLASS5 ADD(); GETCHAR(); SCANTYPE = SLASH
32600           EQ(CLASSN,6)                                   :F(OUT)
32700    CLASS5A REMOVE(); GETNONBLANK()
32800    CLASS5B EQ(CLASSN,6)                                  :F(CLASS5A)
32900           REMOVE(); GETNONBLANK()
33000           EQ(CLASSN,5)                                   :F(CLASS5B)
33100           REMOVE() .                                     : (SCAN)
33200    CLASS6 ADD(); SCANTYPE = STAR                        : (OUT)
33300    CLASS7 ADD(); SCANTYPE = EQUALS                       : (OUT)
33400    CLASS8 ALD(); SCANTYPE = LPAREN                      : (OUT)
33500    CLASS9 ADD(); SCANTYPE = RPAREN                      : (OUT)
33600    CLASS10 ADD(); SCANTYPE = SEMICOLON                  : (OUT)
33700    CLASS11 ADQ(); SCANTYPE = COMMA                     : (OUT)
33800    CLASS12 IMPOSSIBLE
33900    CLASSO REMOVE()
34000      OUTPUT = " /***** ERROR, UNKNOWN CHARACTER **"
34100    +              NEXTCHAR "' ENCOUNTERED; DELETED. *****/" : (SCAN)
34200    *
34300    OUT        UNARY =
34400               UNARY = EQ(SCANTYPE,EQUALS) '1'
34500               UNARY = EQ(SCANTYPE,LPAREN) '1'
34600    +                                                  : (RETURN)
34700    *-----*
34800    GETCHAR CARD LEN(1) . NEXTCHAR                       : F(MORE)
34900             CLASSN = IDENT(NEXTCHAR,' ') 12                : S(RETURN)
3500G         CLASSN = IDENT(NEXTCHAR,' ') 12                 : S(RETURN)
35100         CLASSN = IDENT(NEXTCHAR,',') 11                   : S(RETURN)
35200         CLASSN = REPLACE(NEXTCHAR,
35300   +           'ABCDEFGHIJKLMNOPQRSTUVWXYZ_#0123456789+-/*=( ); ',
35400   +           '00000000000000000000000000000000011111111123456789')
35500         CLASSN = INTEGER(CLASSN) CONVERT(CLASSN,.INTEGER) + 1
35600     +                                                  : S(RETURN)
35700         CLASSN = 0                                          : (RETURN)
35800    MORE      CARD = INPUT                                 : S(GETCHAR)
35900             OUTPUT = ' /***** EOF HIT *****/'
36000             CLASSN = 10
```

36100		NEXTCHAR = ';' .	:(RETURN)
36200	*		
36300	GETNCABLANK	GETCHAR()	
36400		EQ(CLASSN,12)	:F(RETURN)
36500		CARD LEN(1) =	:(GETNONBLANK)
36600	*		
36700	ADD	SCAN = SCAN NEXTCHAR	
36800		CARD LEN(1) =	:(RETURN)
36900	*		
37000	REMOVE	CARD LEN(1) =	:(RETURN)
37100	-EJECT		
37200	END		

SNORTRAN System Routines

```

00100
00200      DECLARE
00300          1  SMEARED_NUMBER,
00400          2  COMPUTED_VALUE  FLOAT BINARY(21),
00500          2  VARIANCE        FLOAT BINARY(21), /*UNITS BETA**(-2T)*/
00600          2  INTDATA,
00700          3  IA_LB           FLOAT BINARY(21), /* INTERVAL */
00800          3  IA_UB           FLOAT BINARY(21); /* ANALYSIS DATA */
00900
01000      DECLARE
01100          1 REG(10) LIKE SMEARED_NUMBER,
01200          BETA_MT  FLOAT BINARY(21) INIT(1.0E-24B), /*16**(-6)*/
01300          HUGE     FLOAT BINARY(21) INIT(1.0E+88B);
01400
01500  SMEARED_ADD: PROC (A,B,C);
01600      /* PERFORMS ADDITION A = B + C FOR SMEARED NUMBERS */
01700      DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,
01800              (T,T1,T2,T3) FLOAT BINARY(53);
01900      T = B.COMPUTED_VALUE;
02000      A.COMPUTED_VALUE = RCUND(T + C.COMPUTED_VALUE);
02100      T2 = MAXSQ(B.INTDATA);
02200      T3 = MAXSQ(C.INTDATA);
02300      CALL INTERVAL(A.INTDATA,1 /*ADD*/, B.INTDATA,C.INTDATA);
02400      T1 = MINSQ(A.INTDATA);
02500      IF ((T2+T3)/T1) > HUGE THEN
02600          DO; A.VARIANCE = HUGE; RETURN; END;
02700      A.VARIANCE = ROUNDUP( T2/T1*B.VARIANCE + T3/T1*C.VARIANCE );
02800      CALL ADD_ROUNDING_ERROR(A);
02900      END SMEARED_ADD;
03000
03100  SMEARED_NEGATE: PROC (A,B);
03200      /* PERFORMS NEGATION A = -B FOR SMEARED NUMBERS */
03300      DECLARE 1 (A,B) LIKE SMEARED_NUMBER,
03400              T FLOAT BINARY(21);
03500      A.COMPUTED_VALUE = -B.COMPUTED_VALUE;
03600      A.VARIANCE = B.VARIANCE;
03700      T = B.IA_LB;
03800      A.IA_LB = -B.IA_UB;
03900      A.IA_UB = -T;
04000      END SMEARED_NEGATE;
04100
04200  SMEARED_SUB: PROC (A,B,C);
04300      /* PERFORMS SUBTRACTION A = B - C FOR SMEARED NUMBERS */
04400      DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,
04500              (T,T1,T2,T3) FLOAT BINARY(53);
04600      T = B.COMPUTED_VALUE;
04700      A.COMPUTED_VALUE = RCUND(T - C.COMPUTED_VALUE);
04800      T2 = MAXSQ(B.INTDATA);
04900      T3 = MAXSQ(C.INTDATA);
05000      CALL INTERVAL(A.INTDATA,2 /*SUB*/, B.INTDATA,C.INTDATA);
05100      T1 = MINSQ(A.INTDATA);
05200      IF ((T2+T3)/T1) > HUGE THEN
05300          DO; A.VARIANCE = HUGE; RETURN; END;
05400      A.VARIANCE = ROUNDUP( T2/T1*B.VARIANCE + T3/T1*C.VARIANCE );
05500      CALL ADD_ROUNDING_ERROR(A);
05600      END SMEARED_SUB;
05700
05800  SMEARED_MULT: PROC (A,B,C);
05900      /* PERFORMS MULTIPLICATION A = B * C FOR SMEARED NUMBERS */
06000      DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,

```

```

06100          T FLOAT BINARY(53);
06200          T = B.COMPUTED_VALUE;
06300          A.COMPUTED_VALUE = RCUND(T * C.COMPUTED_VALUE);
06400          CALL INTERVAL(A.INTDATA,3 /*MULT*/,B.INTDATA,C.INTDATA);
06500          T = B.VARIANCE;
06600          A.VARIANCE = RCUNDUP(T + C.VARIANCE);
06700          CALL ADD_ROUNDING_ERROR(A);
06800          END SMEARED_MULT;
06900
07000 SMEARED_DIV: PROC(A,B,C);
07100          /* PERFORMS DIVISION A = B / C FOR SMEARED NUMBERS */
07200          DECLARE 1 (A,B,C) LIKE SMEARED_NUMBER,
07300                  T FLCAT BINARY(53);
07400          T = B.COMPUTED_VALUE;
07500          A.COMPUTED_VALUE = RCUND(T / C.COMPUTED_VALUE);
07600          CALL INTERVAL(A.INTDATA,4 /*DIV*/,B.INTDATA,C.INTDATA);
07700          T = B.VARIANCE;
07800          A.VARIANCE = RCUNDUP(T + C.VARIANCE);
07900          CALL ADD_ROUNDING_ERROR(A);
08000          END SMEARED_DIV;
08100
08200 SETCONSTANT: PROC(A,C);
08300          /* SETS SMEARED NUMBER A TO THE INITIAL VALUE C */
08400          DECLARE 1 A LIKE SMEARED_NUMBER,
08500                  C FLOAT BINARY(53);
08600          A.COMPUTED_VALUE = RCUND(C);
08700          A.IA_LB = ROUNDDOWN(C);
08800          A.IA_UB = ROUNDUP(C);
08900          A.VARIANCE = 0E08;
09000          IF C /= PRECISION(C,21) THEN
09100              CALL ADD_ROUNDING_ERROR(A);
09200          END SETCONSTANT;
09300
09400 INTERVAL: PROC(A,OP,B,C);
09500          /* PERFORMS INTERVAL ARITHMETIC A = B OP C */
09600          DECLARE 1 INTERVAL,
09700                  2 IA_LB      FLCAT BINARY(21),
09800                  2 IA_UB      FLOAT BINARY(21),
09900                  1 (A,B,C) LIKE INTERVAL,
10000                  (S1,S2,S3,S4) FLOAT BINARY(53),
10100                  IOP(4) LABEL,
10200                  OP          FIXED BINARY(15);
10300          S1 = B.IA_LB; S2 = B.IA_UB; S3 = C.IA_LB; S4 = C.IA_UB;
10400          GO TO ICP(OP);
10500
10600          IOP(1):          /* ADDITION */
10700          A.IA_LB = ROUNDDOWN(S1 + S3);
10800          A.IA_UB = ROUNDUP (S2 + S4);
10900          RETURN;
11000
11100          IOP(2):          /* SUBTRACTION */
11200          A.IA_LB = ROUNDDOWN(S1 - S4);
11300          A.IA_UB = ROUNDUP (S2 - S3);
11400          RETURN;
11500
11600          IOP(3):          /* MULTIPLICATION */
11700          A.IA_LB = MIN( ROUNDDOWN(S1 * S3),
11800                        ROUNDDOWN(S1 * S4),
11900                        ROUNDDOWN(S2 * S3),
12000                        ROUNDDOWN(S2 * S4) );

```



```

12100      A.IA_UB = MAX(   ROUNDUP   (S1 * S3),
12200                        ROUNDUP   (S1 * S4),
12300                        ROUNDUP   (S2 * S3),
12400                        ROUNDUP   (S2 * S4)   );
12500      RETURN;
12600
12700      IOP(4):   /* DIVISION */
12800      A.IA_LB = MIN(   ROUNDDOWN(S1 / S3),
12900                        ROUNDDOWN(S1 / S4),
13000                        ROUNDDOWN(S2 / S3),
13100                        ROUNDDOWN(S2 / S4)   );
13200      A.IA_UB = MAX(   ROUNDUP   (S1 / S3),
13300                        ROUNDUP   (S1 / S4),
13400                        ROUNDUP   (S2 / S3),
13500                        ROUNDUP   (S2 / S4)   );
13600      RETURN;
13700
13800      END INTERVAL;
13900
14000      MAXSQ:  PROC(A) RETURNS(FLOAT BINARY(53));
14100      /* RETURNS MAX VALUE OF A**2 OVER INTERVAL A */
14200      DECLARE 1 A, /* INTERVAL */
14300                2 LB FLOAT BINARY(21),
14400                2 UB FLOAT BINARY(21),
14500                (T1,T2) FLOAT BINARY(53);
14600      T1 = PRECISION(A.LB,53)**2;
14700      T2 = PRECISION(A.UB,53)**2;
14800      T1 = MAX(T1,T2);
14900      RETURN(T1);
15000      END MAXSQ;
15100
15200      MINSQ:  PROC(A) RETURNS(FLOAT BINARY(53));
15300      /* RETURNS MIN VALUE OF A**2 OVER INTERVAL A */
15400      DECLARE 1 A, /* INTERVAL */
15500                2 LB FLOAT BINARY(21),
15600                2 UB FLOAT BINARY(21),
15700                (T1,T2) FLOAT BINARY(53);
15800      IF (A.LB*A.UB <= 0E0B) THEN /* INTERVAL CONTAINS ZERO */
15900          DO; T1 = 1E0B/HUGE; RETURN(T1); END;
16000      T1 = PRECISION(A.LB,53)**2;
16100      T2 = PRECISION(A.UB,53)**2;
16200      T1 = MIN(T1,T2);
16300      RETURN(T1);
16400      END MINSQ;
16500
16600
16700      ADD_ROUNDING_ERROR: PROC(A);
16800      DECLARE 1 A LIKE SNEAKED_NUMBER,
16900                (D,S) FLOAT BINARY(53);
17000      /* COMPUTE      D = MIN(MANTISSA(A))           */
17100      /*              S = 1 / (2D)**2 / 3             */
17200      /*              = COMPUTATION ROUNDOFF ERROR VARIANCE, */
17300      /*              UNITS BETA**(-2T).              */
17400
17500      IF A.VARIANCE >= HUGE | A.COMPUTED_VALUE=0E0B THEN RETURN;
17600      D = MIN(MANTISSA(A.IA_LB),MANTISSA(A.IA_UB));
17700      S = 1E0B / (D**2*12E0);
17800      A.VARIANCE = ROUNDUP(A.VARIANCE + D);
17900      END ADD_ROUNDING_ERROR;
18000

```

```

18100 MANTISSA: PROC(X) RETURNS(FLOAT BINARY(53));
18200 DECLARE (X,Y) FLOAT BINARY(21);
18300 Y = ABS(X);
18400 DO WHILE (Y > 1.06000000000000000000E+08);
18500     I = Y * 1.00000000000000000000E-48;
18600 END;
18700 DO WHILE (Y < 1.00000000000000000000E-48);
18800     Y = Y * 1.00000000000000000000E+48;
18900 END;
19000 RETURN(Y);
19100 END MANTISSA;
19200
19300
19400 SMEARED_PRINT: PROC(AA,A);
19500 DECLARE I A LIKE SMEARED_NUMBER,
19600         AA CHAR(*),
19700         (S,S1,S2,ABS1,ABS2,ABS3,ABS4) FLOAT BIN(53);
19800 S = SQRT(A.VARIANCE);
19900 S1 = -3.625 * S; /* 99.7% CONFIDENCE LIMITS */
20000 S2 = +3.625 * S;
20100 ABS3 = A.COMPUTED_VALUE * (1 + S1*BETA_MT);
20200 ABS4 = A.COMPUTED_VALUE * (1 + S2*BETA_MT);
20300 IF A.COMPUTED_VALUE < 0 THEN
20400     DO; ABS1=ABS3; ABS2=ABS4; ABS3=ABS1; END;
20500
20600 PUT FILE(SYSPRINT) SKIP EDIT
20700
20800 (AA,'COMPUTED VALUE:',A.COMPUTED_VALUE)
20900 (SKIP,COL(7),A,X(8),A,E(25,15,16))
21000
21100 ('3.625*SIGMA BOUNDS ON RELATIVE ERROR:', '(' ,S1, ' ',S2,
21200 ' ) * BETA**(-T)' )
21300 (SKIP,COL(12),A,COL(50),A,2 (E(25,15,16),A))
21400
21500 ('IMPLIED 99.7%-CONFIDENCE BCUNDS:', '(' ,ABS3, ' ',
21600 ABS4, ' )' )
21700 (SKIP,COL(12),A,COL(50),A,2 (E(25,15,16),A))
21800
21900 ('INTERVAL ANALYSIS BCUNDS:', '(' ,A.IA_LB, ' ',A.IA_UB, ')' )
22000 (SKIP,COL(12),A,COL(50),A,2 (E(16,6,7),X(9),A));
22100
22200 END SMEARED_PRINT;
22300
22400
22500 ROUND: PROC(X) RETURNS(FLOAT BINARY(21));
22600 DECLARE (X,HALF) FLOAT BINARY(53),
22700 RX FLOAT BINARY(21),
22800 E BIT(64),
22900 EEXPT BIT(8) DEFINED E POSITION(1),
23000 RBIT BIT(1) DEFINED E POSITION(33),
23100 F BIT(64) INIT('0000000000000000000000000000000000000000000000000000000'),
23200 FEXPT BIT(8) DEFINED F POSITION(1);
23300
23400 E = UNSPEC(X);
23500 IF RBIT THEN DO; /* ROUND */
23600     FEXPT = EEXPT;
23700     UNSPEC(HALF) = F;
23800     RX = X + HALF;
23900 END;
24000 ELSE RX = X; /* TRUNCATE */

```



```

24100      RETURN(RX);
24200      END ROUND;
24300
24400  ROUNDUP: PROC(X) RETURNS(FLOAT BINARY(21));
24500      DECLARE X FLOAT BINARY(53),
24600              RX FLOAT BINARY(21);
24700      IF X > CE08 THEN RETURN(ROUND(X));
24800              ELSE DO; RX=X; RETURN(RX); END;
24900      END ROUNDUP;
25000
25100  ROUNDDOWN: PROC(X) RETURNS(FLOAT BINARY(21));
25200      DECLARE X FLOAT BINARY(53),
25300              RX FLOAT BINARY(21);
25400      IF X < OE08 THEN RETURN(ROUND(X));
25500              ELSE DO; RX=X; RETURN(RX); END;
25600      END ROUNDDOWN;

```

## 7. CONCLUSIONS

The important conclusions of this paper are first, that accumulated roundoff errors tend to cluster towards the center of their bounds like normal densities; second, that high-confidence bounds can be proved; third, that automated statistical analysis is not worthwhile.

The last conclusion is not obvious. Naturally the statistical model will generate better bounds than the worst-case model of Wilkinson. But beyond that no conclusions are obvious. Interval Analysis treats error densities as uniform, although the Central Limit Theorem tells us they are almost normal. We would expect better results from statistical analysis because it not only combines intervals (read "probability densities on intervals"), it also retains information about the resultant error distribution on the new interval. But statistical analysis breaks down because relative errors cannot handle cancellation or numbers near zero, and correlation is almost impossible to monitor practically.

One possible exit from the irrepresentability of zero is to drop relative errors altogether and develop a statistical theory of absolute errors. Observations 2 and 3 of Section 2 provide most of the necessary theory if we now treat numbers as random variables instead of the number's relative errors. Zero then becomes representable. However, it is not clear that an analogue of Theorem 5 in Section 4 can be proved for absolute analysis since multiplication and division no longer involve simple convolution, and besides, the problems of correlated errors cannot be eliminated. Statistical error analysis in general seems of questionable use.

## REFERENCES

1. Feller, W. An Introduction to Probability Theory and Its Applications v. II. NY: John Wiley & Sons, 1966.
2. Gnedenko, B. V. and A. N. Kolmogorov. Limit Distributions for Sums of Independent Random Variables. Reading, Mass.: Addison-Wesley, 1954.
3. Hamming, R. W. "On the Distribution of Numbers" Bell Sys. Tech. J. 49, 8, 1609 (1970).
4. Henrici, P. Elements of Numerical Analysis. NY: John Wiley & Sons, 1964.
5. Huskey, H.D. "On the Precision of a Certain Procedure of Numerical Integration." J. Res. Nat. Bur. Stds., 42, 1, 57 (1949).
6. Kaneko, T. and B. Liu. "On Local Roundoff Errors in Floating-Point Arithmetic." JACM 20, 3, 391 (1973).
7. Lever Brothers. "Vanish Advertisement." Shopper's Weekly Magazine, 17 May 1975.
8. McBride, E. B. Obtaining Generating Functions. NY: Springer-Verlag, 1971.
9. Marasa, J. D. and D. W. Matula. "A Simulative Study of Correlated Error Propagation in Various Finite-Precision Arithmetics." IEEE Trans Comp. C-22, 6, 587 (1973).
10. Papoulis, A. The Fourier Integral and its Applications. NY: McGraw-Hill, 1962.
11. Pennington, R. H. Introductory Computer Methods and Numerical Analysis. Toronto: Macmillan, 1970.
12. Prohorov, Yu. V. and Yu. A. Rozanov. Probability Theory. NY: Springer-Verlag, 1969.
13. Sterbenz, P. H. Floating-Point Computation. Englewood Cliffs, NJ: Prentice-Hall, 1974.
14. Tsao, N. "On the Distribution of Significant Digits and Roundoff Errors." CACM 17, 5, 269 (1974).
15. Uspensky, J. V. Introduction to Mathematical Probability. NY: McGraw-Hill, 1937.
16. Wilkinson, J. H. Rounding Errors in Algebraic Processes. Englewood Cliffs, NJ: Prentice-Hall, 1963.

BIBLIOGRAPHIC DATA SHEET		1. Report No. UIUCDCS-R-76-787	2.	3. Recipient's Accession No.	
Title and Subtitle THE STATISTICAL THEORY OF RELATIVE ERRORS IN FLOATING-POINT COMPUTATION				5. Report Date March 1976	
				6.	
Author(s) Douglass Stott Parker, Jr.				8. Performing Organization Rept. No. UIUCDCS-R-76-787	
Performing Organization Name and Address University of Illinois at Urbana-Champaign Department of Computer Science Urbana, Illinois 61801				10. Project/Task/Work Unit No.	
				11. Contract/Grant No. US NSF DCR73-07980 A02	
2. Sponsoring Organization Name and Address  National Science Foundation Washington, D. C.				13. Type of Report & Period Covered Master's Thesis	
				14.	
5. Supplementary Notes					
6. Abstracts  Experience with interval analysis indicates that the error bounds it generates for most floating-point computations are pessimistic. This is necessary since interval analysis assumes some error in every computation, and is oblivious of both error cancellation and the Central Limit Theorem. It has been questioned whether a possible alternative to this generation of bounds is the use of error probability densities to obtain high-confidence intervals, or statistical bounds, on the computed error. This paper shows that such a probabilistic system for Wilkinsonian relative errors can be derived, but points out numerous crippling drawbacks for this system and statistical error analysis in general.					
7 Key Words and Document Analysis. 17a. Descriptors  Automated error analysis Floating-point computation Interval analysis Roundoff error Statistical error bounds					
7b. Identifiers/Open-Ended Terms					
7c. COSATI Field/Group					
8. Availability Statement  Release Unlimited			19. Security Class (This Report) UNCLASSIFIED		21. No. of Pages 66
			20. Security Class (This Page) UNCLASSIFIED		22. Price

















OCT 17 REC'D





UNIVERSITY OF ILLINOIS-URBANA

510.84 IL6R no. C002 no.782-787(1976

Design of Irredundant MOS Networks : a p



3 0112 088402588